

# A Novel Lossless Encoding Algorithm for Data Compression

Abdelghani Tbakhi (✉ [atbakhi@khcc.jo](mailto:atbakhi@khcc.jo))

King Hussein Cancer Center

Anas Al-okaily

King Hussein Cancer Center <https://orcid.org/0000-0002-1188-8266>

---

## Article

**Keywords:** data compression, encoding algorithm, DNA data

**Posted Date:** May 24th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-499681/v1>

**License:**   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

# **A Novel Lossless Encoding Algorithm for Data Compression**

Anas Al-okaily<sup>1</sup>, Abdelghani Tbakhi<sup>1</sup>

<sup>1</sup>*Department of Cell Therapy & Applied Genomics, King Hussein Cancer Center, Amman, 11941, Jordan.*

**Data compression is a challenging and increasingly important problem. As the amount of data generated daily continues to increase, efficient transmission and storage has never been more critical. In this study, a novel encoding algorithm is proposed, motivated by the compression of DNA data and associated characteristics. The proposed algorithm follows a divide-and-conquer approach by scanning the whole genome, classifying subsequences based on similarity patterns, and binning similar subsequences together. The data are then compressed in each bin independently. This approach is different than the currently known approaches: entropy, dictionary, predictive, or transform based methods. Proof-of-concept performance was evaluated using a benchmark dataset with seventeen genomes ranging in size from kilobytes to gigabytes. The results showed considerable improvement in the compression of each genome, preserving several megabytes compared with state-of-art tools. Moreover, the algorithm can be applied to the compression of other data types include mainly text, numbers, images, audio, and video which are being generated daily and unprecedentedly in massive volumes.**

## 20 Introduction

21 The importance of data compression, a fundamental problem in computer science, information the-  
22 ory, and coding theory, continues to increase as global data quantities expand rapidly. The primary  
23 goal of compression is to reduce the size of data for subsequent storage or transmission. There are  
24 two common types of compression algorithms: lossless and lossy. Lossless algorithms guarantee  
25 exact restoration of the original data, while lossy algorithms do not. Such losses are caused, for  
26 instance, by the exclusion of unnecessary information, such as metadata in video or audio that will  
27 not be observed by users. This study focuses specifically on lossless compression. Data exist in  
28 different formats including text, numbers, images, audio, and video. Several coding algorithms  
29 and corresponding variants have been developed for textual data, the primary focus of this pa-  
30 per. This includes the Huffman <sup>1</sup>, Shannon <sup>2</sup>, Shannon-Fano <sup>3</sup>, Shannon-Fano-Elias <sup>4</sup>, Lempel-Ziv  
31 (LZ77)<sup>5</sup>, the Burrows-Wheeler transform <sup>6</sup>, and Tunstall <sup>7</sup> algorithms. The Huffman algorithm  
32 includes several variants: minimum variance Huffman, canonical Huffman, length-limited Huff-  
33 man, non-binary Huffman, adaptive Huffman, Faller-Gallager-Knuth (an adaptive Huffman) <sup>8</sup>, and  
34 Vitter (an adaptive Huffman) <sup>9</sup>. The LZ algorithm also includes several variants, such as LZ78  
35 <sup>10</sup>, Lempel-Ziv-Welch (LZW) <sup>11</sup>, Lempel-Ziv-Stac (LZS) <sup>12</sup>, Lempel-Ziv-Oberhumer (LZO) <sup>13</sup>,  
36 Lempel-Ziv-Storer-Szymanski (LZSS) <sup>14</sup>, Lempel-Ziv-Ross-Williams (LZRW) <sup>15</sup>, and the Lem-  
37 pel-Ziv-Markov chain algorithm (LZMA) <sup>16</sup>. Additional techniques involve arithmetic encoding  
38 <sup>17</sup>, range encoding <sup>18</sup>, move-to-front encoding (also referred as symbol ranking encoding) <sup>19,20</sup>, run-  
39 length encoding <sup>21</sup>, delta encoding, unary encoding, context tree weighting encoding <sup>22</sup>, prediction  
40 by partial matching <sup>23</sup>, context mixing <sup>24</sup>, asymmetric numeral systems (also called asymmetric

41 binary coding)<sup>25</sup>, length index preserving transform<sup>26</sup>, and dynamic Markov encoding<sup>27</sup>.

42 Common tools used to implement one or more encoding algorithms and compress textual  
43 data are listed and described in Table S1. The tools namely are: bcm, blzpack, brotli, bsc, bzip2,  
44 cmix, compress, freeze, gzip, hook, Huffman-codec, lizard, lrzip, lz4, lzb, lzfse, lzip, lzop, lzturbo,  
45 Nakamichi, pbzip2, pigz, ppm, qzip, rans static, rzip, snzip, srnk, xz, zlib, zip, zpipe, and zstd.  
46 The performance of these tools is promising, but further improvements are still possible and nec-  
47 essary, due to the high volume of data that being generated worldwide. The results produced by  
48 these tools were compared with those of the proposed technique as part of the study. Compres-  
49 sion algorithms can be classified based on the methodology used in the algorithm, such as entropy,  
50 dictionary, predictive, and transform based methods. These methods have been described exten-  
51 sively in several recent studies<sup>28–32</sup>, however, a brief description for each method is provided in  
52 the Supplementary Information.

53 Genomics (DNA/RNA) data is a type of textual information with several unique charac-  
54 teristics. First, the alphabet consists only of A, C, G, and T characters representing the four  
55 nucleotides: adenine, cytosine, guanine, and thymine, respectively. Second, DNA data contain  
56 repeat sequences and palindromes. Third, the size of genomics data can be very large, relative to  
57 most media files. The human genome, for instance, consists of more than three billion nucleotides  
58 (specifically 3,272,116,950 bases<sup>33</sup> requiring more than three gigabytes of storage). As such,  
59 sequencing genomic data (especially for humans) is currently being performed for research and  
60 diagnostic purposes in daily basis. In addition, this sequencing is typically conducted with high

61 depth (30-100x) to sequence the same region several times in order to make reading the DNA re-  
62 gions more accurate. This generates massive quantities of data on a daily basis. For example, the  
63 number of bases sequenced from December 1982 through February 2021 was 12,270,717,209,612  
64 <sup>34</sup>. Several algorithms have been developed to compress these data, which can be divided into  
65 vertical and horizontal techniques <sup>35</sup>. Vertical mode algorithms utilize a reference genome/source,  
66 while horizontal mode algorithms are reference-free. The algorithm proposed in this study is a  
67 reference-free (horizontal) model.

68 Genomic data are stored in different formats, including FATSA <sup>36</sup>, FASTQ <sup>37</sup>, and SAM <sup>38</sup>,  
69 with FASTA being the most common and also the primary focus of this paper. Several comparative  
70 studies for compressing FASTA files have been published in recent years <sup>39-43</sup>. Genomic sequences  
71 typically consist of millions or billions of sequenced reads, with lengths in the hundreds, stored  
72 with the quality of each base in a primarily FASTQ format. A common DNA data processing  
73 strategy involves aligning the sequenced reads with a reference genome. The output is the reads  
74 themselves, with their base qualities and alignment results for each read, stored in a SAM format.  
75 Surveys of compression tools for SAM and FASTQ data are available in the literature <sup>40,44</sup>.

76 The small alphabet found in DNA data simplifies the compression process. However, the  
77 problem remains challenging due to the discrete, uniform distribution (frequencies) of bases in  
78 DNA data. Efficient compression relies mainly on repetitiveness in the data and encoding as few  
79 characters/words as possible, since encoding more characters costs more bits-per-character. If the  
80 characters are uniformly distributed in the text, their repetitions will also be distributed uniformly

81 and encoding only a fraction of them (to decrease the bits-per-character) will lead to low compres-  
82 sion outcomes. The application of Huffman encoding, for instance, produces 2-bit assignments for  
83 each base. The algorithm will then produce an inefficient/suboptimal compression result that does  
84 not utilize repetitions found in the DNA data. Motivated by this issue, we introduce a novel loss-  
85 less, reference-free encoding algorithm that is applicable to any data type (i.e., numbers, images,  
86 audio, and video), in addition to genomic or textual data.

## 87 **Methods**

88 The following observations can be inferred from a careful analysis of DNA. First, many regional  
89 (local) sub-sequences (assume a length of 100 bases) contain non-uniform or skewed distributions.  
90 Second, similar sub-sequences (which provide better compression results if encoded together) are  
91 often distant from each other. This distance is typically longer than the length of sliding windows  
92 (usually in kilobases/kilobytes) commonly used in algorithms such as LZ, context weighting tree,  
93 predictive partial matching, or dynamic Markov compression. Even if these sequences are located  
94 within the same sliding window, they are often sufficiently distant from each other, which leads  
95 to inefficient compression and encoding. These two observations were the motivation behind the  
96 design of the following encoding algorithm.

97 **Compression algorithm** Given a string  $S$  of length  $s$ , containing characters from a fixed alphabet  
98 of length  $\Sigma$ , a window of length  $w$ , and a label sequence  $L$  initialized as empty, the proposed  
99 algorithm can be expressed as follows.

- 100 1. Scan  $S$  with non-overlapping windows of length  $w$ . In each window:
- 101 (a) Using a classification method, classify the sequence in the window based on its content.
- 102 Create a unique label  $l$  that denotes this classification.
- 103 (b) Create, if not any, a bin labeled  $l$ . Then, concatenate sequence of the window with the
- 104 sequence in this bin.
- 105 (c) Concatenate  $l$  with  $L$ . The purpose of  $L$  is to store the queue/order of the sequences'
- 106 labels. This order is needed during the decompression process to restore the sequences
- 107 (distributed in the bins) to their original placements in the input data.
- 108 The time cost for this step is  $O(\frac{sc}{w})$  where  $c$  is the cost of classifying each sequence of length
- 109  $w$ . If  $c$  is less than or equal to  $w$ , the cost of this step is  $O(s)$ .
- 110 2. Once the scanning of  $S$  is completed, encode all collected labels in  $L$  using Huffman en-
- 111 coding. The time cost for this step is  $O(B \log B)$  where  $B$  is the number of labels. These
- 112 labels can be set dynamically or designed in advance, and they can be set by the user or the
- 113 implementer.
- 114 3. 3. Compress  $L$  using the encoding schema from the previous step. The time cost of this step
- 115 is  $O(\frac{s}{w})$ .
- 116 4. Compress the sequence in each bin using an algorithm suitable for the content and/or the
- 117 label of the bin.
- 118 5. Pack or archive, with or without further compression, all resultant compressed files (bins and
- 119  $L$ ) together into a single file.

120 Note the value of  $w$  can be fixed or variable. If  $w$  is a variable, the window is extended  
121 character by character until the label of the sequence in the window matches one of the bin labels.  
122 The length of the window must then be added to  $L$  (after the label string for each window) or  
123 placed in a separate stream that stores lengths and their order. Lastly, lengths can be encoded  
124 using a universal code for integers (such as Levenshtein coding <sup>45</sup>, Elias coding <sup>46</sup> (delta, gamma,  
125 or omega), exponential-Golomb code <sup>47</sup>, Fibonacci code <sup>48</sup>, Stout code <sup>49</sup>) or using a suitable  
126 encoding algorithm such as unary, binary, or Huffman.

127 **Decompression algorithm** Decompression is the inverse of compression and can be described as  
128 follows.

- 129 1. If all bins and  $L$  were compressed/archived into a single file, decompress/unarchive them.
- 130 2. Decompress each bins file and  $L$  with the compression tool that was used to compress them.
- 131 3. Initialize a counter for each bin.
- 132 4. Sequentially read labels from  $L$ . At each read, extract a sequence of length  $w$  from the bin  
133 with label equal to the read label. Extraction of the sequence must start from the position  
134 equal to the value of the counter associated with that bin. Increment *counter* by the value of  
135  $w$  plus 1 and then output the extracted sequence to the output stream.

136 The time and memory cost of decompression is linear.

137 This algorithm can be applied not only to DNA or textual data, but to archiving and other data

138 types as well (e.g., numbers, images, audio, and video). Sub-binning or nested-binning processes  
139 can also be applied to each bin.

140 This design facilitates organizing and sorting the input data using a divide-and-conquer  
141 method by creating bins for similar data and encodes/compresses data in the same bin that are  
142 better of compressed together, to achieve better compression results with a minor increase in time  
143 costs. In the case of more random/divergent data, which is common, this is done to avoid relying  
144 on a single encoding or compression technique (as in entropy methods), being dependent on the  
145 randomness of previous sequences (as in prediction methods), requiring construction of dictionar-  
146 ies dependent on the randomness of previous sequences (as in dictionary methods), or potentially  
147 introducing inefficient transformation due to the randomness of the data (as in transform methods).  
148 In contrast, the proposed algorithm divides the data into groups of similar segments, regardless of  
149 their position in the original sequence, which decreases the randomness and contributes in orga-  
150 nizing the input data to ultimately handling the compression process more efficiently.

151 Note that the continued application of sub-binning processes will eventually reduce the ran-  
152 domness/divergence of the data and improve the compression results, by obtaining data that are  
153 optimal or suboptimal for compression. These processes will require additional time costs, but  
154 these costs will still be practical at low sub-binning depth and feasible at high sub-binning depths,  
155 especially for small data or the compression of large data for archiving. Therefore, sub-binning  
156 will eventually provide more control, organization, and possibly a deterministic solution to en-  
157 coding and compression problems. Further analysis and investigations are also provided in the

158 Supplementary Information.

159 This encoding algorithm is named in honor of Pramod Srivastava (Professor in the Depart-  
160 ment of Immunology, University of Connecticut School of Medicine) who served as an advisor to  
161 the first author. As such, it is named the Okaily-Srivastava-Tbakhi (*OST*) algorithm.

162 **OST-DNA** The first implementation of the OST algorithm (OST-DNA) accepts DNA data as in-  
163 put. Bin labels are computed using a Huffman tree encoding strategy. For example, if the Huffman  
164 tree for a subsequence produces the following encoding schema: G:0, A:10, T:110, and C:111,  
165 then the label will be GATC\_1233 label (1 indicates G is encoded using 1 bit, A using 2 bits, and  
166 so on). The number of bits used to encode a character gives a general indication of its frequency  
167 compared to the other characters. The number of bins can be reduced by decreasing the label  
168 length as follows. To produce a label length of 1, we used the first base of the Huffman tree and  
169 its bit length. As such, the above Huffman encoding schema will be represented by G\_1. If the bin  
170 label length is 2, then the label will be GA\_12. This clearly decreases the number of labels, but at  
171 the cost of decreasing the similarity among sequences in the same bin therefore their compression.

172 As the windows do not overlap, each base in  $S$  will be read in  $O(1)$  time. The cost of  
173 constructing a Huffman tree for a subsequence is then  $O(\Sigma \log \Sigma)$ , requiring the construction of  
174  $O(\frac{s}{w})$  Huffman trees for all subsequences in  $S$ . The total cost hence is  $O(\frac{s \Sigma \log \Sigma}{w})$ . In order to allow  
175 for the acquisition of non-uniform distributions for the characters in  $\Sigma$  (the pigeonhole principal),  
176 the value of  $w$  must be larger than that of  $\Sigma \log \Sigma$ , noting that  $\Sigma$  is a constant. As such, the total  
177 cost the compression process of OST-DNA is  $O(s)$ .

178 Since the value of  $w$  is fixed in this version of OST-DNA, Huffman trees are constructed once  
179 for each window sequence. In the case of a variable  $w$  where the window will be extended until  
180 the sequence label matches a bin label, it is not efficient to calculate Huffman trees for the entire  
181 sequence at every extension, hence adaptive Huffman trees can be applied instead. Generally, the  
182 compressed bin files and  $L$  can be collected into a single file using an archiver which could perform  
183 further compression. However, this process was omitted in this study to demonstrate the efficiency  
184 of the OST algorithm without any further compression may be produced by the archiving process.

## 185 **Results**

186 We implemented OST-DNA using the python language. We used the same dataset applied to  
187 another benchmark <sup>39</sup> in order to test and compare the compression results from OST-DNA with  
188 the tools listed in Table S1. The dataset consists of seventeen genomes, as shown in Tables S2 and  
189 S3. The main purpose of this implementation and benchmarking is to provide a proof-of-concept  
190 version of the proposed algorithm for academic use.

191 The following preprocessing steps were applied to each tested genome. All new lines, header  
192 lines, lowercase bases, and bases not identified as A, C, G, T, or N, were removed. This produced a  
193 one-line sequence for each genome, containing only the four capitalized DNA bases and the letter  
194 “N”. The python script used to perform these preprocessing steps and the size of each genome,  
195 before and after applying the script, are provided in Table S3. The size of one-line genomes ranged  
196 from 50 KB to 13 GB, with a total size of 16,773.88 Megabytes.

197 Compression ratio was the primary metric used for evaluating the performance of the pro-  
198 posed algorithm. It is equal to the size of the compressed file divided by the size of the uncom-  
199 pressed (original) file. The original files in this study are the one-line genome files. Other metrics  
200 include compression time (seconds), decompression time (seconds), compression speed (the size  
201 of the uncompressed file in MB divided by the compression time in seconds-MB/s), and the de-  
202 compression speed (the size of the uncompressed file in MB divided by the decompression time in  
203 seconds-MB/s).

204 Each tool in Table S1 was applied to each one-line genome. The compression and decom-  
205 pression commands used to run each tool are provided in Table S5. The cumulative compression  
206 results (for all one-line genomes) are provided in Table 1, while the compression results for each  
207 one-line genome are listed in Extended Data Table 1. The most efficient tools in terms of com-  
208 pression ratio were lrzip (saved 14,317.40 MB out of 16,773.88 MB), brotli (13,958.42 MB), lzip  
209 (13,916.39 MB), xz (13,915.92 MB), bsc (13,391.09 MB), and bcm (13,314.83 MB). In addition,  
210 comparing the results of the commonly used tools (bzip2 and gzip) indicated bzip2 was better,  
211 saving 12,601.12 MB.

212 Seven versions of OST-DNA were implemented. In each version, one of the seven most  
213 efficient tools (bcm, brotli, bsc, lrzip, lzip, xz, and bzip2) is used to compress the bins generated  
214 by the OST-DNA algorithm. The same command used by each tool to compress one-line genomes  
215 was used to compress the bins. Each of these seven versions were run over each window lengths  
216 of 25, 50, 100, 125, 150, 250, 500, 750, 1,000, 2,500, 5,000, and 10,000 and across label lengths

217 of 1, 2, 3, 4, and 5. The best cumulative compression performance achieved by each OST-DNA  
218 version is shown in Table 2.

219 A comparison of the results produced by each OST-DNA tool (i.e., bcm, brotli, bsc, bzip2,  
220 lrzip, lzip, and xz) indicated OST-DNA-bcm saved an additional 77.38 MB compared to bcm, OST-  
221 DNA-brotli: 140.41 MB, OST-DNA-bsc: 66.79 MB, OST-DNA-bzip2: 34.83 MB, OST-DNA-  
222 lrzip: 12.05 MB, OST-DNA-lzip: 38.34 MB, and OST-DNA-xz: 41.65 MB. This demonstrates the  
223 proposed algorithm can improve compression results compared to the individual corresponding  
224 tools.

225 The best tool in terms of compression ratio was lrzip, yet OST-DNA-lrzip saved an addi-  
226 tional 12.05 MB more than lrzip. In terms of compression time, bsc was the fastest tool. OST-  
227 DNA-bsc could save an additional 66.79 MB more than bsc with a practical increase in the com-  
228 pression/decompression times (hence corresponding decrease in compression and decompression  
229 speeds). These increases are a result of the time needed for classifying and binning sequences  
230 during compression, as well as the need to collect and restore the original genome during decom-  
231 pression. However, they can be decreased significantly as follows. First, the OST-DNA script was  
232 not optimized for implementation but was intended in this study to provide proof of concept. Ad-  
233 ditional improvements to the script can reduce both the compression and decompression times by  
234 increasing the corresponding speeds. In addition, parallel processing, which could further reduce  
235 run-time, was not applied during the binning, compression, or decompression steps of OST-DNA.  
236 Finally, fewer bins would lead to faster sequence labeling and longer window lengths could speed

237 up both compression and decompression, with a trade-off in the compression ratio.

238 The compression results for OST-DNA using each of the seven tools for each one-line  
239 genome were also better than the results using the corresponding standalone tool. This can be  
240 found by comparing the compression results using each OST-DNA version with each window  
241 length and each label length for the one-line genomes, as shown in Extended Data Table 2. The  
242 compression results for the standalone tools are provided in Table 1.

243 By analyzing the compression results for all OST-DNA versions, using different sequence  
244 label lengths and classification methods (i.e., Huffman tree encoding schema), we found the most  
245 efficient results correlated with a window length of 250 to 1,000 bases. This is reasonable, as  
246 lengths shorter or longer than this will yield a uniform distribution of bases in the sequence. How-  
247 ever, dynamic window lengths can be more practical and feasible given the additional costs for  
248 encoding the lengths. We found efficient label lengths to be 2 and 4. This is reasonable as in-  
249 creased label lengths produce more bins and more similarities among sequences in the same bin.  
250 Compression is more efficient when sequences in a bin are more similar. Additional classification  
251 methods can be used to label the bins and improve compression further. Extended Data Table 3  
252 shows compression results for each version of OST-DNA for each window and each label length  
253 cumulatively applied to all one-line genomes. Further analysis at the bin level, rather than the  
254 genome level, is provided in Extended Data Table 4.

255 Compression results produced by applying each OST-DNA version to each bin, using the  
256 same window and label lengths but with different genomes, were considerably different (see Ex-

257 tended Data Table 5). This was not the case for bins produced using the same label length and  
258 same genome, but with different window lengths (see Extended Data Table 6). This means that  
259 sequences with the same label but from different genomes differed significantly (even though their  
260 labels were the same). This observation suggests the need to find a set of labels or labeling steps  
261 that could compress sequences from any source (genome) with similar efficiency, to improve the  
262 compression results further. In other words, sequences that share a label from this set would be  
263 compressed at a similar rate, regardless of the source (genome) from which they were derived.  
264 This set of labels could also be used better archival of multiple genomes.

265 The current version of OST-DNA compresses each bin using a specific tool. However, this  
266 is not optimal. Finding a tool that optimally compresses each bin, or a novel algorithm that is  
267 customized for efficient compression based on the bin content or label, could further improve the  
268 overall performance.

## 269 **References**

- 270 1. Huffman, D. A. *et al.* A method for the construction of minimum-redundancy codes. *Pro-*  
272 *ceedings of the IRE* **40**, 1098–1101 (1952).
- 273 2. Shannon, C. E. A mathematical theory of communication. *ACM SIGMOBILE mobile com-*  
274 *puting and communications review* **5**, 3–55 (2001).
- 275 3. Fano, R. M. *The transmission of information* (Massachusetts Institute of Technology, Research  
276 Laboratory of Electronics . . . , 1949).

- 277 4. Cover, T. M. *Elements of information theory* (John Wiley & Sons, 1999).
- 278 5. Ziv, J. & Lempel, A. A universal algorithm for sequential data compression. *IEEE Transac-*  
279 *tions on information theory* **23**, 337–343 (1977).
- 280 6. Burrows, M. & Wheeler, D. J. A block-sorting lossless data compression algorithm. *Citeseer*  
281 (1994).
- 282 7. Tunstall, B. P. *Synthesis of noiseless compression codes*. Ph.D. thesis, Georgia Institute of  
283 Technology (1967).
- 284 8. Knuth, D. E. Dynamic huffman coding. *Journal of algorithms* **6**, 163–180 (1985).
- 285 9. Vitter, J. S. Design and analysis of dynamic huffman codes. *Journal of the ACM (JACM)* **34**,  
286 825–845 (1987).
- 287 10. Ziv, J. & Lempel, A. Compression of individual sequences via variable-rate coding. *IEEE*  
288 *transactions on Information Theory* **24**, 530–536 (1978).
- 289 11. Welch, T. A. A technique for high-performance data compression. *Computer* 8–19 (1984).
- 290 12. Friend, R. C. Transport Layer Security (TLS) Protocol Compression Using Lempel-Ziv-Stac  
291 (LZS). RFC 3943 (2004). URL <https://rfc-editor.org/rfc/rfc3943.txt>.
- 292 13. Oberhumer, M. Lzo-a real-time data compression library. [http://www.oberhumer.](http://www.oberhumer.com/opensource/lzo/)  
293 [com/opensource/lzo/](http://www.oberhumer.com/opensource/lzo/) (2008).
- 294 14. Storer, J. A. & Szymanski, T. G. Data compression via textual substitution. *Journal of the*  
295 *ACM (JACM)* **29**, 928–951 (1982).

- 296 15. Williams, R. N. An extremely fast ziv-lempel data compression algorithm. In [1991] *Pro-*  
297 *ceedings. Data Compression Conference*, 362–371 (IEEE, 1991).
- 298 16. Ranganathan, N. & Henriques, S. High-speed vlsi designs for lempel-ziv-based data compres-  
299 sion. *IEEE Transactions on Circuits and Systems II: Analog and Digital Signal Processing*  
300 **40**, 96–106 (1993).
- 301 17. Langdon, G. G. An introduction to arithmetic coding. *IBM Journal of Research and Develop-*  
302 *ment* **28**, 135–149 (1984).
- 303 18. Martín, G. Range encoding: an algorithm for removing redundancy from a digitised message.  
304 In *Video and Data Recording Conference, Southampton, 1979*, 24–27 (1979).
- 305 19. Ryabko, B. Y. Data compression by means of a “book stack”. *Problemy Peredachi Informatsii*  
306 **16**, 16–21 (1980).
- 307 20. Bentley, J. L., Sleator, D. D., Tarjan, R. E. & Wei, V. K. A locally adaptive data compression  
308 scheme. *Communications of the ACM* **29**, 320–330 (1986).
- 309 21. Capon, J. A probabilistic model for run-length coding of pictures. *IRE Transactions on*  
310 *Information Theory* **5**, 157–163 (1959).
- 311 22. Willems, F. M., Shtarkov, Y. M. & Tjalkens, T. J. The context-tree weighting method: basic  
312 properties. *IEEE transactions on information theory* **41**, 653–664 (1995).
- 313 23. Cleary, J. & Witten, I. Data compression using adaptive coding and partial string matching.  
314 *IEEE transactions on Communications* **32**, 396–402 (1984).

- 315 24. Mahoney, M. V. Adaptive weighing of context models for lossless data compression. Tech.  
316 Rep., Florida Tech (2005).
- 317 25. Duda, J. Asymmetric numeral systems: entropy coding combining speed of huffman coding  
318 with compression rate of arithmetic coding. *arXiv preprint arXiv:1311.2540* (2013).
- 319 26. Awan, F. S. & Mukherjee, A. Lipt: A lossless text transform to improve compression. In  
320 *Proceedings International Conference on Information Technology: Coding and Computing*,  
321 452–460 (IEEE, 2001).
- 322 27. Cormack, G. V. & Horspool, R. N. S. Data compression using dynamic markov modelling.  
323 *The Computer Journal* **30**, 541–550 (1987).
- 324 28. Gopinath, A. & Ravisankar, M. Comparison of lossless data compression techniques. In *2020*  
325 *International Conference on Inventive Computation Technologies (ICICT)*, 628–633 (IEEE,  
326 2020).
- 327 29. Kavitha, P. A survey on lossless and lossy data compression methods. *International Journal*  
328 *of Computer Science & Engineering Technology (IJCSET)* **7** (2016).
- 329 30. squash. <https://quixdb.github.io/squash-benchmark/>.
- 330 31. Uthayakumar, J., Vengattaraman, T. & Dhavachelvan, P. A survey on data compression tech-  
331 niques: From the perspective of data quality, coding schemes, data type and applications.  
332 *Journal of King Saud University-Computer and Information Sciences* (2018).

- 333 32. Kodituwakku, S. & Amarasinghe, U. Comparison of lossless data compression algorithms for  
334 text data. *Indian journal of computer science and engineering* **1**, 416–425 (2010).
- 335 33. humangenomesize. [https://www.ncbi.nlm.nih.gov/grc/human/data?asm=](https://www.ncbi.nlm.nih.gov/grc/human/data?asm=GRCh38.p13)  
336 [GRCh38.p13](https://www.ncbi.nlm.nih.gov/grc/human/data?asm=GRCh38.p13).
- 337 34. genbank. <https://www.ncbi.nlm.nih.gov/genbank/statistics/>.
- 338 35. Grumbach, S. & Tahi, F. A new challenge for compression algorithms: genetic sequences.  
339 *Information Processing & Management* **30**, 875–886 (1994).
- 340 36. Lipman, D. J. & Pearson, W. R. Rapid and sensitive protein similarity searches. *Science* **227**,  
341 1435–1441 (1985).
- 342 37. Cock, P. J., Fields, C. J., Goto, N., Heuer, M. L. & Rice, P. M. The sanger fastq file format for  
343 sequences with quality scores, and the solexa/illumina fastq variants. *Nucleic acids research*  
344 **38**, 1767–1771 (2010).
- 345 38. Li, H. *et al.* The sequence alignment/map format and samtools. *Bioinformatics* **25**, 2078–2079  
346 (2009).
- 347 39. Kryukov, K., Ueda, M. T., Nakagawa, S. & Imanishi, T. Sequence compression benchmark  
348 (scb) database—a comprehensive evaluation of reference-free compressors for fasta-formatted  
349 sequences. *GigaScience* **9**, giaa072 (2020).
- 350 40. Hosseini, M., Pratas, D. & Pinho, A. J. A survey on data compression methods for biological  
351 sequences. *Information* **7**, 56 (2016).

- 352 41. Mansouri, D., Yuan, X. & Saidani, A. A new lossless dna compression algorithm based on a  
353 single-block encoding scheme. *Algorithms* **13**, 99 (2020).
- 354 42. Bakr, N. S., Sharawi, A. A. *et al.* Dna lossless compression algorithms. *American Journal of*  
355 *Bioinformatics Research* **3**, 72–81 (2013).
- 356 43. Jahaan, A., Ravi, T. & Panneer Arokiaraj, S. A comparative study and survey on existing dna  
357 compression techniques. *International Journal of Advanced Research in Computer Science* **8**  
358 (2017).
- 359 44. Bonfield, J. K. & Mahoney, M. V. Compression of fastq and sam format sequencing data. *PloS*  
360 *one* **8**, e59190 (2013).
- 361 45. Levenshtein. [http://www.compression.ru/download/articles/int/](http://www.compression.ru/download/articles/int/levenstein_1968_on_the_redundancy_and_delay.pdf)  
362 [levenstein\\_1968\\_on\\_the\\_redundancy\\_and\\_delay.pdf](http://www.compression.ru/download/articles/int/levenstein_1968_on_the_redundancy_and_delay.pdf), author=Vladimir  
363 Levenshtein, year=1968.
- 364 46. Elias, P. Universal codeword sets and representations of the integers. *IEEE transactions on*  
365 *information theory* **21**, 194–203 (1975).
- 366 47. Salomon, D. *Data compression: the complete reference* (Springer Science & Business Media,  
367 2004).
- 368 48. Fraenkel, A. S. & Kleinb, S. T. Robust universal complete codes for transmission and com-  
369 pression. *Discrete Applied Mathematics* **64**, 31–55 (1996).

- 370 49. Stout, Q. Improved prefix encodings of the natural numbers (corresp.). *IEEE Transactions on*  
371 *Information Theory* **26**, 607–609 (1980).

Table 1: Compression performance for each common tool cumulatively over all tested genomes.

Tool	CP Ratio (%)	Saved Space (MB)	CP Time (seconds)	DP Time (seconds)	CP Speed (MB/s)	DP Speed (MB/s)
bcm	20.6217	13,314.83	4,071	3,728	4.1203	0.9279
blzpack	37.2279	10,529.31	227	149	73.8937	41.9098
brotli	16.7848	13,958.42	62,659	103	0.2677	27.3346
bsc	20.1670	13,391.09	2,369	68	7.0806	49.7469
bzip2	24.8765	12,601.12	2,474	1,169	6.7801	3.5695
compress	25.3977	12,513.70	443	168	37.8643	25.3582
freeze	27.4616	12,167.51	6,078	233	2.7598	19.7698
gzip	26.9031	12,261.19	4,211	171	3.9833	26.3900
hook	21.3395	13,194.41	7,803	8,074	2.1497	0.4433
Huffman-codec	27.4015	12,177.59	1,152	401	14.5607	11.4621
lizard	34.8186	10,933.45	11,449	41	1.4651	142.4494
lrzip	14.6446	14,317.40	21,589	378	0.7770	6.4986
lz4	52.7757	7,921.34	161	73	104.1856	121.2677
lzfse	29.3101	11,857.45	1,118	130	15.0035	37.8187
lzip	17.0353	13,916.39	20,079	304	0.8354	9.3996
lzop	47.6212	8,785.96	152	106	110.3545	75.3578
LzTurbo	28.4807	11,996.56	157	46	106.8400	103.8548
pbzip2	24.9054	12,596.28	2,458	1,193	6.8242	3.5018
pigz	26.9356	12,255.74	4,284	109	3.9155	41.4508
ppm	23.8049	12,780.87	5,020	6,314	3.3414	0.6324
qzip	41.9873	9,730.98	1,556	126	10.7801	55.8960
rans	24.0431	12,740.93	144	91	116.4853	44.3182
rzip	24.9315	12,591.89	2,515	1,279	6.6695	3.2697
snzip	45.5241	9,137.73	159	73	105.4961	104.6048
srank	41.2318	9,857.70	794	778	21.1258	8.8897
xz	17.0381	13,915.92	18,666	266	0.8986	10.7442
zip	26.9031	12,261.19	4,165	173	4.0273	26.0849
zlib	26.9187	12,258.57	4,278	117	3.9210	38.5924
zpipe	26.9187	12,258.57	4,283	106	3.9164	42.5973
zstd	26.7482	12,287.18	251	75	66.8282	59.8227

CP is abbreviation for compression and DP for decompression. The size of all genomes is 16,773.88 MB. The tools cmix, lzb, and Nakamichi could not compress large genomes in reasonable time so their cumulative performance could not be presented.

Table 2: Compression performance for best window-length and label-length of each of the seven OST-DNA versions cumulatively over all tested genomes.

Tool	Window Length	Label Length	CP Ratio (%)	Saved Space (MB)	CP Time (seconds)	DP Time (seconds)	CP Speed (MB/s)	DP Speed (MB/s)
<b>bcm</b>	-	-	20.6217	13,314.83	4,071	3,728	4.1203	0.9279
<b>OST-DNA-bcm</b>	250	4	20.1603	13,392.21	12,026	4,388	1.3948	0.7707
<b>brotli</b>	-	-	16.7848	13,958.42	62,659	103	0.2677	27.3346
<b>OST-DNA-brotli</b>	750	4	15.9477	14,098.83	72,315	318	0.2320	8.4121
<b>bsc</b>	-	-	20.1670	13,391.09	2,369	68	7.0806	49.7469
<b>OST-DNA-bsc</b>	250	5	19.7688	13,457.88	10,355	2,160	1.6199	1.5352
<b>bzip2</b>	-	-	24.8765	12,601.12	2,474	1,169	6.7801	3.5695
<b>OST-DNA-bzip2</b>	250	2	24.6689	12,635.95	10,132	1,230	1.6555	3.3642
<b>lrzip</b>	-	-	14.6446	14,317.40	21,589	378	0.7770	6.4986
<b>OST-DNA-lrzip</b>	1,000	1	14.5728	14,329.45	31,911	432	0.5256	5.6584
<b>lzip</b>	-	-	17.0353	13,916.39	20,079	304	0.8354	9.3996
<b>OST-DNA-lzip</b>	750	2	16.8067	13,954.74	30,691	472	0.5465	5.9728
<b>xz</b>	-	-	17.0381	13,915.92	18,666	266	0.8986	10.7442
<b>OST-DNA-xz</b>	750	2	16.7898	13,957.57	29,098	393	0.5765	7.1662

CP is abbreviation for compression and DP for decompression. The size of all genomes is 16,773.88 MB.

372 **Competing Interests** None declared.

373 **Correspondence** Correspondence and requests for materials should be addressed to Anas Al-okaily (email:

374 AA.12682@khcc.jo).

375 **Authors Contribution** Both authors contributed equally to this work and hence both names should be

376 considered first name.

377 **Computer code** Source code of the seven OST-DNA tools are available at <https://github.com/>

378 aalokaily/OST.

379 **Supplementary Information** is available for this paper

380 **Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints)

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Supplementaryinformation.pdf](#)
- [ExtendedDataTable1.xlsx](#)
- [ExtendedDataTable2.xlsx](#)
- [ExtendedDataTable3.xlsx](#)
- [ExtendedDataTable4.xlsx](#)
- [ExtendedDataTable5.xlsx](#)
- [ExtendedDataTable6.xlsx](#)