

# Coordinate-based mapping of tabular data enables fast and scalable queries

Stephen R. Piccolo<sup>1,\*</sup>, Zachary E. Ence<sup>1</sup>, Kimball Hill<sup>1</sup>, PJ Tatlow<sup>1</sup>, Brandon J. Fry<sup>1</sup>, Jonathan B. Dayton<sup>1</sup>

1 - Department of Biology, Brigham Young University, Provo, UT, 84602, USA

\* - To whom correspondence should be addressed.

## Abstract

**Motivation:** Biologists commonly store data in tabular form with observations as rows, attributes as columns, and measurements as values. Due to advances in high-throughput technologies, the sizes of tabular datasets are increasing. Some datasets contain millions of rows or columns. To work effectively with such data, researchers must be able to efficiently extract subsets of the data (using filters to select specific rows and retrieving specific columns). However, existing methodologies for querying tabular data do not scale adequately to large datasets or require specialized tools for processing. We sought a methodology that would overcome these challenges and that could be applied to an existing, text-based format.

**Results:** In a systematic benchmark, we tested 10 techniques for querying simulated, tabular datasets. These techniques included a delimiter-splitting method, the Python *pandas* module, regular expressions, object serialization, the *awk* utility, and string-based indexing. We found that storing the data in fixed-width formats provided excellent performance for extracting data subsets. Because columns have the same width on every row, we could pre-calculate column and row coordinates and quickly extract relevant data from the files. Memory mapping led to additional performance gains. A limitation of fixed-width files is the increased storage requirement of buffer characters. Compression algorithms help to mitigate this limitation at a cost of reduced query speeds. Lastly, we used this methodology to transpose tabular files that were hundreds of gigabytes in size, without creating temporary files. We propose coordinate-based, fixed-width storage as a fast, scalable methodology for querying tabular biological data.

*Contact:* stephen\_piccolo@byu.edu

## 29 Introduction

30 Biologists often generate data suitable for representation in an attribute-value system<sup>1</sup>, also known  
31 as an information system<sup>2</sup>, simple frame<sup>3</sup>, object-predicate table<sup>4</sup>, or flat file. In this representation,  
32 an *object* might be a biological organism, an *attribute* might be a characteristic of that organism,  
33 and a *value* might be a datum for that object and attribute. For example, a researcher might observe  
34 200 cancer patients (objects) and collect transcriptomic measurements for 20,000 genes (attributes);  
35 each *value* would indicate the relative number of transcripts present in tumor cells for each  
36 patient/gene combination<sup>5</sup>. In this example, the data values would have been summarized  
37 previously using preprocessing tools, such as a reference aligner and a transcript-quantification  
38 algorithm<sup>6-9</sup>. For convenience and compactness, researchers typically store attribute-value data in  
39 2-dimensional, tabular formats. Commonly, in such tables, each row contains data for a given  
40 object, and each column contains data for a given attribute<sup>10</sup>; but in some cases, the table is  
41 transposed (objects as columns, attributes as rows). Researchers use tabular data to perform  
42 analytical tasks, such as executing statistical analyses, producing graphics, and further summarizing  
43 the data.

44 Across the subfields of biology, researchers store a considerable proportion of tabular data in  
45 plain-text formats. This approach coincides with the Unix and “Pragmatic Programming”  
46 philosophies<sup>11-13</sup>, which advocate for storing data and sharing data among computer programs as  
47 plain text. Keeping data as text has many advantages. Plain text is readable to humans.  
48 Sophisticated text editors are freely available for all major operating systems. A wide range of tools  
49 exist for generating, manipulating, parsing, and compressing text files; these include  
50 long-established command-line tools developed by the Unix community. In addition, scripting  
51 languages like Python<sup>14</sup> and R<sup>15</sup> provide libraries for analyzing text-based tabular data; these  
52 libraries are used broadly within the biology community and elsewhere<sup>16</sup>. Storing data as plain text  
53 does have drawbacks relative to binary formats. Text files may be larger than binary files, and it  
54 may be more computationally intensive to parse a text file than a binary file. Consequently, a  
55 multitude of techniques for storing tabular data in binary (non-text) formats has been developed.  
56 For example, researchers use Microsoft Excel for data exploration and analysis<sup>17,18</sup>; relational  
57 databases provide a formalized methodology to query tabular data<sup>19</sup>; so-called NoSQL databases  
58 provide alternative methodologies for structuring and querying data, including attribute-value  
59 systems<sup>20</sup>; the Hierarchical Data Format (HDF5) is often used for tabular data, including in biology  
60 research<sup>21,22</sup>. Additionally, in recent years, distributed architectures for large-scale data storage and  
61 processing have seen wide use; these technologies include Apache Hadoop and Apache Spark<sup>23,24</sup>.  
62 Despite these advances, the humble plain-text file continues to play a critical role in biology  
63 research due to its simplicity, flexibility, familiarity, and portability.

64 In our own research studying molecular profiles of tumors—and via collaborations with other  
65 scientists—we have frequently encountered a need to *select* and *project* tabular data. In  
66 relational-algebra terms<sup>19</sup>, selection refers to the process of identifying rows that match some  
67 criteria; projection refers to the process of retrieving specific columns. For example, in a study of  
68 genomic and transcriptomic profiles of human breast tumors, a researcher might wish to select only  
69 patients diagnosed before the age of 40 and who harbor a mutation in *BRCA1*, a gene known to  
70 effect double-stranded DNA repair by homologous recombination<sup>25</sup>. Having identified this patient  
71 subset, the researcher might wish to retrieve (project) transcriptomic data for genes that interact  
72 with *BRCA1*. In repositories like The Cancer Genome Atlas (TCGA), The International Genome  
73 Consortium (ICGC), and Gene Expression Omnibus (GEO)<sup>26-28</sup>, genomic and transcriptomic

74 data—and their corresponding annotations—are stored in tabular text files. The ways that objects,  
75 attributes, and values are oriented within these files differ across these and other repositories<sup>29–31</sup>,  
76 but values are commonly oriented in rows and columns and are separated by tab characters, comma  
77 characters, or some other delimiter.

78 Typically, to parse such data, researchers write custom scripts or use software packages that  
79 facilitate parsing<sup>32–34</sup>. To perform selection, the code must extract all values from the column(s) to  
80 be used as filtering criteria. If data values are delimited by tab characters, for example, the code  
81 must identify the positions of tab characters and extract data at the relevant positions for each row.  
82 However, because data values may vary in length, the positions of tab characters may differ for each  
83 row, and these positions must be reidentified for each row, thus slowing execution. After identifying  
84 rows that match the selection criteria, the researcher may then wish to project the data. When  
85 parsing a tab-delimited file, the code must again identify positions of tab characters *for each row*  
86 and extract values at the relevant positions. In this methodology, the code parses the data row by  
87 row, thus minimizing memory consumption. Alternatively, the entire file could be parsed into an  
88 in-memory data structure; this methodology may increase the efficiency of selection and projection,  
89 but many datasets are too large to fit in memory. Additionally, if a researcher wishes to use only a  
90 few columns for selection or projection, it is inefficient to read the entire file into memory. Hybrid  
91 solutions exist, such as the *pandas* module for Python<sup>35</sup>. However, in our experience, it has been  
92 difficult to find a solution that strikes a satisfactory balance between speed and memory usage. This  
93 challenge has become more acute as data sizes have increased. For example, a re-quantification of  
94 RNA expression data from TCGA contains data for 11,373 tumors across 199,169 transcripts<sup>36</sup>.  
95 Phase I of the Library of Integrated Network-based Cellular Signatures (LINCS) yielded data for  
96 more than 1.3 million experiments, including transcriptomic data for 12,300 genes (after  
97 imputation) and annotations for each experiment<sup>37</sup>. Recently, the UK Biobank posted genotypic,  
98 phenotypic, and health-related measurements for approximately 500,000 individuals<sup>38</sup>. Some of  
99 these files are multiple terabytes in size. These trends are true in other fields as well, including  
100 proteomics, remote sensing, and imaging<sup>39–42</sup>.

101 In evaluating methodologies that could handle such data, we envisioned scenarios in which data  
102 files are created once and then queried many times. Public repositories like TCGA, LINCS, and UK  
103 Biobank cater to these scenarios; after the data have been prepared, they are stored on web servers,  
104 enabling researchers to download and query the data. Because the files are written only once, it is  
105 less important to optimize speeds for writing the files, and it is unnecessary to support concurrent  
106 writing by multiple agents. In contrast, it is highly preferable that researchers can query the data  
107 quickly and flexibly. With this context in mind, we sought a solution that would meet the following  
108 criteria:

- 109 • Handle datasets larger than what can fit into memory on modern personal computers.
- 110 • Handle attributes of different types (categorical, ordinal, numeric, etc.).
- 111 • Support selection based on data in *any* column.
- 112 • Store the data in a portable format that can be transported across systems without custom  
113 tools or specialized expertise.
- 114 • Store the data in a space-efficient manner (while preferring fast speeds over reduced storage).
- 115 • Not be specific to any particular type of biological data (e.g., genomic, transcriptomic,  
116 ecological).
- 117 • Can be created, indexed, and queried in a non-proprietary<sup>43</sup>, programming-language agnostic,  
118 and platform-independent manner.

- 119 • Can transpose rows and columns without reading *all* the data into memory and without  
120 creating temporary files.
- 121 • Can represent missing values explicitly.

122 In our quest to identify a solution that would address these criteria, we considered a variety of  
123 binary-based solutions. These included relational databases, NoSQL databases, HDF5, and the  
124 Apache Parquet format<sup>44</sup>. With each solution, we faced limitations. For example, the SQLite  
125 relational database has a limit of 32,767 columns<sup>45</sup>. NoSQL databases provide many options for  
126 structuring the data, but we failed to identify an approach that would provide adequate query speeds  
127 and storage sizes. The HDF5 format is designed primarily for numerical data, whereas we sought  
128 the ability to handle other data types as well. As a columnar storage solution, Parquet was efficient  
129 at projection; however, it was ill-suited to selection. Ultimately, we focused on text-based solutions,  
130 performing a benchmark analysis of 10 different techniques for parsing tabular data. As described  
131 below, we chose one of these techniques and refined it further. We found that this technique  
132 addresses each of the above criteria yet is human readable, fast, and scalable.

## 133 Methods

134 In an initial round of benchmarks, we used Python scripts to generate tabular text files in which  
135 10% of the columns contained categorical values (randomly generated, 2-digit alphabetical  
136 sequences) and 90% of the columns contained numerical values (ranging between 0.0 and 1.0).  
137 First, we used these scripts to generate relatively small files, containing 100 columns and 1000  
138 rows. After verifying functionality, we generated two types of large file that represent dimensions  
139 that will be increasingly seen in biological research: 1) “tall” files containing 1 million rows and  
140 1,000 columns, and 2) “wide” files containing 1,000 rows and 1 million columns. Each of these  
141 files contained a total of 1 billion data points (approximately 10 GB in size). For each set of  
142 dimensions, we saved the data in four different formats:

- 143 • *tsv*. We separated each value on each row with tabs (tab-separated-value format).
- 144 • *msgpack*. We used the MessagePack format<sup>46</sup> to serialize each row of data as a list object.
- 145 • *flags*. In an attempt to make it faster to access elements at a given column index, we specified  
146 the index of each element within each row of data and embedded these indices within the file,  
147 prior to each datum.
- 148 • *fwf*. The width of each column corresponded to the data value with the largest number of  
149 characters in that column (fixed-width format). We also added a buffer character between  
150 columns.

151 In this phase, we evaluated 10 techniques for projecting the data. Different techniques used  
152 different versions of the input data (see below). We coded each technique to select the first column  
153 and every hundredth column thereafter. Each script saved the selected columns to a tab-delimited  
154 text file. We then used a script to verify that the output was correct.

- 155 • *delimiter-split*. We used TSV files as input, split each line on tabs, and extracted values at the  
156 specified indices.
- 157 • *pandas*. We applied the *read\_csv* function from the Python *pandas* module to the TSV files.
- 158 • *reg-ex-quant*. We used regular expressions to quantify tab characters that preceded each  
159 specified index and then extracted those values using capturing groups.

- 160 • *reg-ex-tab*. We used regular expressions to map non-capturing groups to indices that should  
161 be ignored and capturing groups to indices that should be extracted.
- 162 • *msgpack*. We deserialized each row from the MessagePack serialized files and extracted  
163 values at the specified indices.
- 164 • *flags*. For the flag files, we identified the position of each specified flag and then extracted  
165 characters after it until another flag was reached.
- 166 • *awk*. We applied the *awk* command-line utility to the TSV files<sup>47</sup>. This Unix-based tool  
167 provides extensive support for parsing text files.
- 168 • *gawk*. This is another variation on *awk*.
- 169 • *nawk*. This is yet another variation on *awk*.
- 170 • *fixed-width*. We used the header line in the file to identify the starting and ending positions of  
171 each column and then used string indexing to extract values at the specified indices.

172 Aside from the *awk*-based solutions, we used Python code. In addition, for each Python solution,  
173 we implemented a memory-mapping version of the code. Memory mapping supports the ability to  
174 randomly access locations within a file; accordingly, in some cases, we could extract specific  
175 portions of the file without needing to iterate sequentially through the file or read every character  
176 into memory.

177 During the second phase of this study, we developed a modified version of the fixed-width format  
178 (*fwf2*). First, we calculated the position and width of each column and stored these values in an  
179 index file that we could also memory map. Second, we calculated the full length of the first  
180 line—all lines should have the same length—and stored this length in a second index file. These  
181 changes enabled us to quickly calculate row and column coordinates when querying the data. Lastly,  
182 we removed the extraneous buffer characters between the columns. This reduced file sizes; however,  
183 we retained a nonessential newline character at the end of each line to make the files more readable.

184 All of our code, along with a bash script to execute the benchmarks, can be found at  
185 [https://github.com/srp33/Tabular\\_File\\_Benchmark](https://github.com/srp33/Tabular_File_Benchmark). The same repository contains an R Markdown  
186 file that includes the code we used to create figures for this paper. We used R version 3.5.1 and the  
187 *ggplot2*, *readr*, *dplyr*, and *cowplot* packages for the figures<sup>48–51</sup>. For the benchmarks, we used  
188 Python (version 3.6.7) and the following external Python modules: *msgpack* (0.5.6), *numpy*  
189 (1.15.2), *pandas* (0.23.4), and *snappy* (0.5.3). All benchmark tests were executed on a 64-bit  
190 processor running Ubuntu Linux (18.04) with the following hardware specifications:

- 191 • 4.5 GHz Xeon® W-2155 (3.3 up to 4.5 GHz – 10 Cores - 20 Threads - 2666 MHz)
- 192 • 256 GB Quad Channel DDR4 random access memory at 2666 MHz (8× 32GB)
- 193 • 250 GB NVMe PCIe M.2 solid-state drive (SSD) for the operating system
- 194 • 3.8 TB NVMe 3.84TB U.2 Mixed Use SSD for data storage

## 195 Results

196 First, we evaluated methods for projecting tabular text files that contained 1 billion data points.  
197 These files had either a “tall” or “wide” orientation. The tall files simulate scenarios in which  
198 researchers collect 1,000 data points for 1 million patients (or other object type). The wide files  
199 simulate scenarios in which researchers collect 1 million data points for 1,000 patients. As  
200 high-throughput data-generation technologies advance and as researchers combine individual  
201 datasets into aggregate ones, such scenarios will be increasingly common.



202 Commonly, biology researchers store data in tab- or comma-delimited files and parse such files  
203 using the *delimiter-split* method. Thus, we considered the performance of this approach to be a  
204 baseline. For the tall files, our scripts extracted every hundredth column in 21.72 and 17.76 seconds  
205 with and without memory mapping, respectively. All but two of the competing methods  
206 outperformed this approach (Figure 1). In contrast, on wide files, the performance slowed  
207 considerably for all methods. The baseline method extracted every hundredth column in 31.38 and  
208 27.44 seconds. The *pandas* method and both regular-expression methods performed worse than the  
209 baseline, and their performance was dramatically worse than it was on the tall files. The poor  
210 performance of *pandas* is perhaps surprising, given the package’s popularity among data  
211 scientists<sup>52</sup>.

212 The fixed-width method performed best overall on the tall file, projecting the data in only 5.34  
213 seconds with memory mapping; however, its performance was mediocre on the wide file. We  
214 hypothesized that a few adjustments to the file format and our algorithmic approach might improve  
215 the performance substantially (see Methods). In addition, we implemented a chunking scheme in  
216 which we parsed 1000 rows of input data at a time before writing to the output file. After these  
217 adjustments, we projected every hundredth row from the tall file in 3.70 seconds. For the wide file,  
218 we projected the data in 3.43 seconds, only 10% the duration of the original approach. Given these  
219 results, we focused on this method and evaluated its performance further.

220 We wanted a method that would excel at projection *and* selection. Therefore, we performed  
221 selection on data from one column with categorical data and one column with numerical data. The  
222 categorical values were 2-character sequences of letters; arbitrarily, we searched for values that  
223 started with “A” or ended with “Z”. The numerical filter searched the specified column in the  
224 remaining rows for values greater than or equal to 0.1. These criteria yielded approximately 6.9%  
225 of the rows. Lastly, we projected every hundredth column. This process took 1.04 seconds for the  
226 tall file and 0.28 seconds for the wide file.

227 When performing the initial benchmarks, we stored the data in four tabular formats (see Methods).  
228 The *flags* and *fixed-width* files were larger than the other formats, especially for the wide files (see  
229 Figure 2). To enable indexing, these formats require extra text within the files. We considered ways  
230 to reduce this extra storage requirement while still supporting fast query times. We tested four  
231 compression algorithms: *gzip*, *bzip2*, *lzma*, and *snappy*. We compressed the text files line by line.  
232 After compression, the lengths of the lines varied, so we saved the starting position of each row to a  
233 serialized dictionary. Compression times differed considerably across the methods (Figure 3); as its  
234 name implies, *snappy* was extremely fast. In contrast, *snappy*-compressed files were  
235 approximately twice as large as files compressed using the other algorithms. Most importantly for  
236 this study, select-and-project speeds were dramatically faster for *snappy*-compressed files than for  
237 any of the other algorithms (Figure 3). However, these speeds are 20-50 times slower than we  
238 attained using non-compressed data. We needed to decompress each full line before we could  
239 evaluate the selection criteria or perform projection. Accordingly, individuals who consider using  
240 this methodology must consider the substantial tradeoff between speed and space requirements.

241 Next we tested the scalability of the (non-compressed) fixed-width approach. We simulated a  
242 scenario in which a researcher might wish to store genotypes for a large number of individuals. We  
243 generated text files that contained simulated genotypes with 10, 50, 100, 500, 1000, 5000, 10000,  
244 50000, 100000 or 500000 rows and columns, respectively. In these files, we represented genetic  
245 loci as rows and organisms as columns and used a pair of nucleotide characters (A, C, G, T) to  
246 represent each genotype. We tested our ability to select and project the data by extracting genotypes  
247 for the intersection of 10 random rows and 10 random columns. Query speeds were identical (0.03

248 seconds) for all file sizes except the largest (0.06 seconds). The largest file (500000 rows by 500000  
249 columns) had a total of 250 billion data points and was 465 GB in size (Figure 4A). Although  
250 extracting 10 rows and 10 columns does not reflect a real-world scenario, it illustrates the promise  
251 of performing these operations quickly on extremely large files.

252 As a final test, we transposed the simulated genotype files. To our knowledge, no tool exists for  
253 transposing tabular files without reading the full dataset into memory or writing temporary files.  
254 Our fixed-width approach successfully transposed each of the simulated genotype files, including  
255 the largest, without writing temporary files. Time and memory usage did increase as file sizes  
256 increased, but in a nonlinear fashion (Figure 4B).

## 257 Discussion

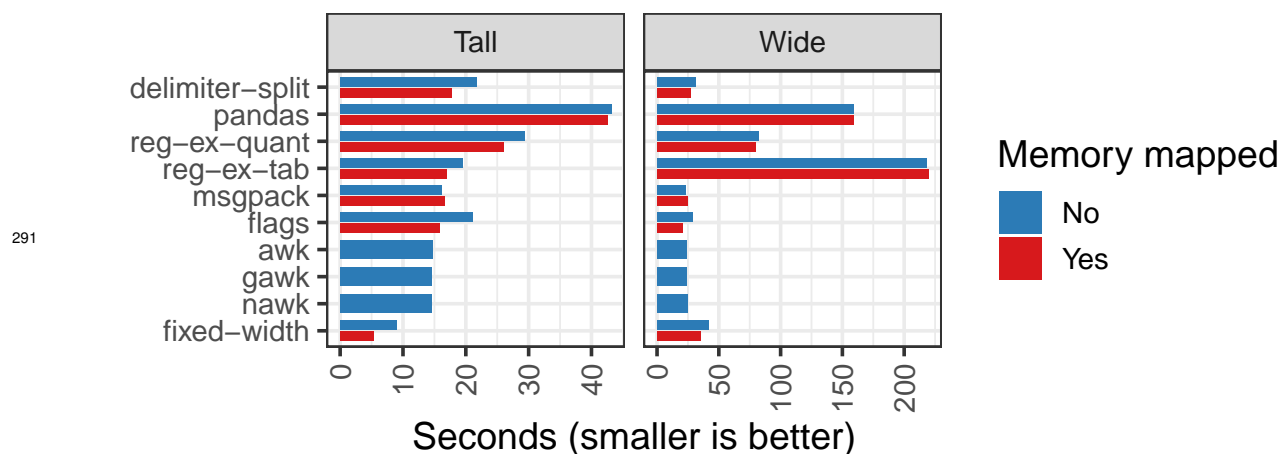
258 Our goal was to evaluate techniques for storing and querying tabular data. Such data are used  
259 widely in biology research. We sought to identify a methodology that would reduce query times  
260 and overcome limitations of existing approaches. In addition, we sought a solution that would  
261 provide the flexibility and readability of plain text plus some benefits of binary formats. Our results  
262 suggest that fixed-width formatting, memory mapping, and coordinate-based indices can meet these  
263 needs for many research scenarios, especially those that prioritize fast reading over fast writing and  
264 that can accept larger file sizes as a compromise for faster querying. In this study, we simulated  
265 data for which all values in a given column are the same width; but in practice, data values in a  
266 given column often vary in length. In these scenarios, extra buffer characters are needed. More  
267 research is necessary to evaluate how much these buffer characters would increase file sizes in  
268 practice, but we predict that query speeds will be impacted only minimally. One possibility for  
269 mitigating the effects of larger file sizes is to compress the whole file using a standard compression  
270 scheme (e.g., gzip) before it is placed on a web server; accordingly, distributing the file would  
271 require less disk space on the server and less network bandwidth during file transfer; the researcher  
272 could then decompress the file locally before querying it.

273 The results of time-based benchmarks must be interpreted with caution because execution times  
274 vary from one computer system to another. Additionally, we performed these benchmarks on  
275 hardware whose performance exceeds that of many computer systems used currently for biology  
276 research. However, we are confident in our conclusions about the *relative* performances of the  
277 methods we evaluated.

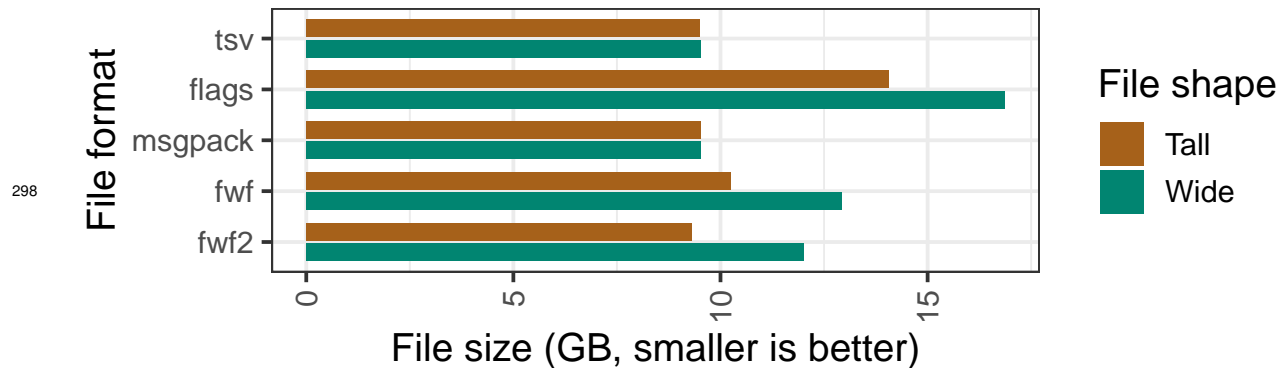
278 This study describes a proof of concept rather than a production-ready tool. One limitation of our  
279 current approach is that data types are not stored explicitly for each attribute; these must be inferred  
280 from the data. However, we believe our methodology's performance in these benchmarks merits  
281 further development.

282 The simplicity of our methodology is one of its strong points. It should be possible to implement  
283 this approach in any programming language and operating system that support reading and writing  
284 text files as well as memory mapping. We do not intend this to be "yet another file format" for  
285 bioinformaticians to deal with; rather, we describe it as a methodology for extending an existing  
286 format. In addition, our approach could facilitate translation among formats and data orientations.  
287 Further methodological refinements are possible, potentially including more sophisticated  
288 algorithms for identifying column widths and compressing the data. We invite collaborations with  
289 others in the research community.

## 290 Figures



**Figure 1: Execution speed for 10 methods of parsing tabular text files.** We evaluated techniques for projecting every hundredth column from tabular text files. The techniques varied in the ways that text files were structured and how they were parsed (see Methods). *awk*, *gawk*, and *nawk* did not support memory mapping. Tall files consisted of 1 million rows and 1,000 columns; “Wide” files had 1,000 rows and 1 million columns. Each file included a mixture of categorical (10%) and numerical (90%) attributes. Note that the x-axis scales differ for the two panels.

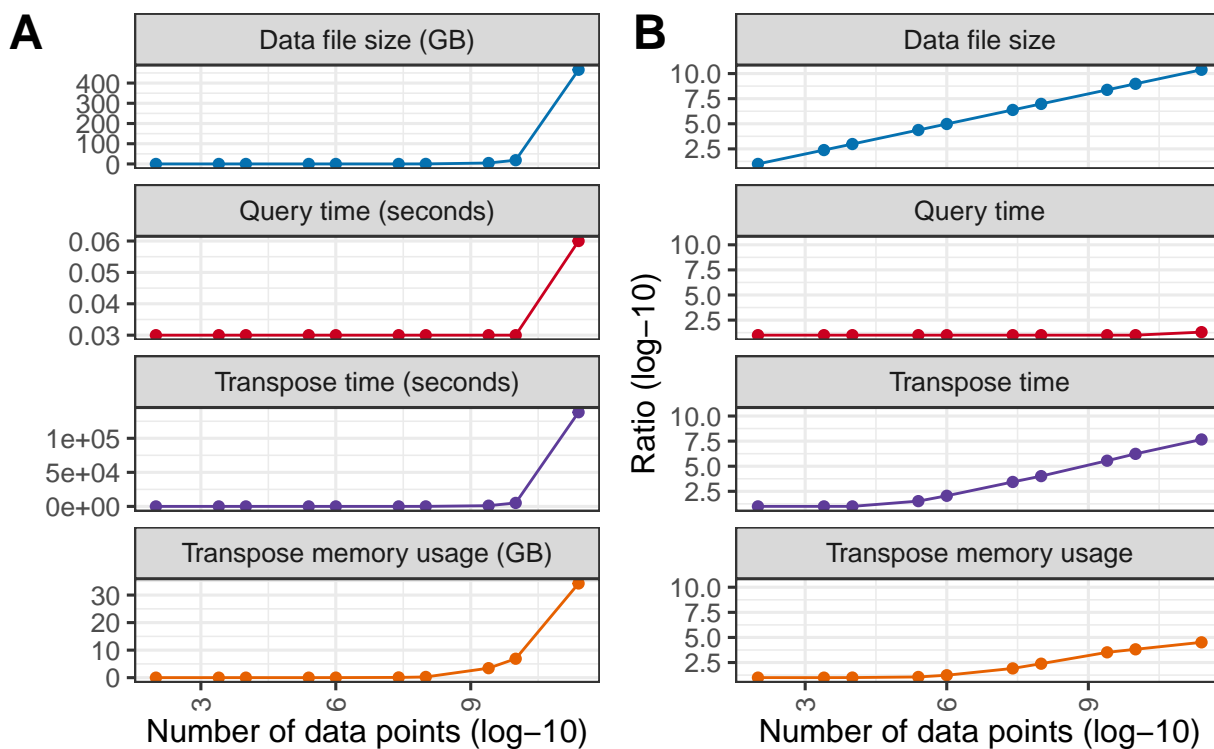


**Figure 2: Sizes of simulated data files used in the initial benchmarks.** *tsv* = Tab-separated values. *flags* = A “flag” before each value indicates the column index of each value. *msgpack* = Each row was serialized as a Python list into MessagePack format. *fwf* = Fixed-width format. *fwf2* = Modified fixed-width format.





303 **Figure 3: Results of compression benchmarks.** In an attempt to reduce file sizes, we compressed the  
 304 fixed-width (fwf2) files line by line and stored the starting position of each line to enable faster file traversal.  
 305 The first panel shows how long it took to compress the files. The second panel indicates file sizes after  
 306 compression (the original files were approximately 10 GB in size). The third panel illustrates how long it  
 307 took to select and project data from the compressed files. The gzip and bzip2 algorithms support a parameter  
 308 to alter the level of compression; we used levels 1 and 9, which are indicated in parentheses. Although the  
 309 *snappy* compression algorithm was much faster than the other algorithms, these speeds were 20-50X slower  
 310 than without compression.



312 **Figure 4: Results of simulated genotype benchmarks.** We simulated genotypes for cohorts and genomes  
313 of increasing size. The x-axes indicate the total number of simulated genotypes. In **A**, the y-axes indicate  
314 *absolute* performance for each metric. In **B**, the y-axes indicate performance *relative* to what was observed  
315 for the minimum number of genotypes. File size increased linearly, whereas the other metrics, especially  
316 query time, increased at a slower rate.

## 317 References

- 318 1. Ziarko, W. & Shan, N. A method for computing all maximally general rules in attribute-value systems.  
319 *Computational Intelligence* **12**, 223–234 (1996).
- 320 2. Pawlak, Z. Information systems theoretical foundations. *Information Systems* **6**, 205–218 (1981).
- 321 3. Barsalou, B., Lawrence W & Hale, C. R. Components of conceptual representation: From feature lists to  
322 recursive frames. in *Categories and Concepts: Theoretical Views and Inductive Data Analysis* 97–144  
323 (Academic Press, 1993).
- 324 4. Watanabe, S. *Pattern Recognition: Human and Mechanical*. (John Wiley & Sons, Inc., 1985).
- 325 5. Mortazavi, A., Williams, B. A., McCue, K., Schaeffer, L. & Wold, B. Mapping and quantifying  
326 mammalian transcriptomes by RNA-Seq. *Nature methods* **5**, 621–8 (2008).
- 327 6. Liao, Y., Smyth, G. K. & Shi, W. The Subread aligner: Fast, accurate and scalable read mapping by  
328 seed-and-vote. *Nucleic acids research* **41**, e108 (2013).
- 329 7. Liao, Y., Smyth, G. K. & Shi, W. FeatureCounts: An efficient general purpose program for assigning  
330 sequence reads to genomic features. *Bioinformatics* **30**, 923–930 (2014).
- 331 8. Patro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon provides fast and bias-aware  
332 quantification of transcript expression. *Nature Methods* **14**, 417–419 (2017).
- 333 9. Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic RNA-seq quantification.  
334 *Nature Biotechnology* **34**, 525–527 (2016).
- 335 10. Wickham, H. Tidy Data. *Journal of Statistical Software* **59**, (2014).
- 336 11. McIlroy, M. D., Pinson, E. N. & Tague, B. A. UNIX Time-Sharing System: Foreword. *Bell System*  
337 *Technical Journal* **57**, 1899–1904 (1978).
- 338 12. Raymond, E. S. *The Art of UNIX Programming*. (Addison-Wesley Professional, 2003).
- 339 13. Hunt, A. & Thomas, D. *The Pragmatic Programmer: From Journeyman to Master*. (Addison-Wesley  
340 Professional, 1999).
- 341 14. Van Rossum, G. & others. Python Programming Language. in *USENIX Annual Technical Conference* **41**,  
342 36 (2007).
- 343 15. R Core Team. *R: A Language and Environment for Statistical Computing*. (R Foundation for Statistical  
344 Computing, 2019).
- 345 16. Stein, L. Creating a bioinformatics nation. *Nature* **417**, 119 (2002).
- 346 17. Ogata, Y. *et al.* KAGIANA: An Excel-Based Tool for Retrieving Summary Information on Arabidopsis  
347 Genes. *Plant and Cell Physiology* **50**, 173–177 (2009).

- 348 18. Simon, R. *et al.* Analysis of Gene Expression Data Using BRB-Array Tools. *Cancer Informatics* **3**,  
349 117693510700300022 (2007).
- 350 19. Codd, E. F. A Relational Model of Data for Large Shared Data Banks. *Commun. ACM* **13**, 377–387  
351 (1970).
- 352 20. Leavitt, N. Will NoSQL Databases Live Up to Their Promise? *Computer* **43**, 12–14 (2010).
- 353 21. Folk, M., Heber, G., Koziol, Q., Pourmal, E. & Robinson, D. An Overview of the HDF5 Technology  
354 Suite and Its Applications. in *Proceedings of the EDBT/ICDT 2011 Workshop on Array Databases* 36–47  
355 (ACM, 2011). doi:[10.1145/1966895.1966900](https://doi.org/10.1145/1966895.1966900)
- 356 22. Enache, O. M. *et al.* The GCTx format and cmap { }Py, R, M, J{ } packages: Resources for optimized  
357 storage and integrated traversal of annotated dense matrices. *Bioinformatics* (2018).
- 358 23. Apache Software Foundation. Apache Hadoop. <https://hadoop.apache.org>
- 359 24. Apache Software Foundation. Apache Spark - Unified Analytics Engine for Big Data.  
360 <https://spark.apache.org>
- 361 25. Moynahan, M. E., Chiu, J. W., Koller, B. H. & Jasin, M. Brca1 Controls Homology-Directed DNA  
362 Repair. *Molecular Cell* **4**, 511–518 (1999).
- 363 26. Koboldt, D. C. *et al.* Comprehensive molecular portraits of human breast tumours. *Nature* **490**, 61–70  
364 (2012).
- 365 27. International Cancer Genome Consortium *et al.* International network of cancer genome projects. *Nature*  
366 **464**, 993–998 (2010).
- 367 28. Barrett, T. *et al.* NCBI GEO: Archive for functional genomics data sets 10 years on. *Nucleic acids*  
368 *research* **39**, D1005–10 (2011).
- 369 29. Rayner, T. F. *et al.* A simple spreadsheet-based, MIAME-supportive format for microarray data:  
370 MAGE-TAB. *BMC Bioinformatics* **7**, 489 (2006).
- 371 30. Reich, M. *et al.* GenePattern 2.0. *Nature genetics* **38**, 500–1 (2006).
- 372 31. Giardine, B. *et al.* Galaxy: A platform for interactive large-scale genome analysis. *Genome research* **15**,  
373 1451–5 (2005).
- 374 32. Davis, S. & Meltzer, P. S. GEOquery: A bridge between the Gene Expression Omnibus (GEO) and  
375 BioConductor. *Bioinformatics (Oxford, England)* **23**, 1846–7 (2007).
- 376 33. Morgan, M. & Davis, S. R. GenomicDataCommons: A Bioconductor Interface to the NCI Genomic Data  
377 Commons. *bioRxiv* (2017). doi:[10.1101/117200](https://doi.org/10.1101/117200)
- 378 34. Colaprico, A. *et al.* TCGAAbiolinks: An R/Bioconductor package for integrative analysis of TCGA data.  
379 *Nucleic Acids Research* **44**, e71–e71 (2016).
- 380 35. McKinney, W. Data Structures for Statistical Computing in Python. in *Proceedings of the 9th Python in*  
381 *Science Conference* (eds. van der Walt, S. & Millman, J.) 51–56 (2010).
- 382 36. Tatlow, P. J. & Piccolo, S. R. A cloud-based workflow to quantify transcript-expression levels in public  
383 cancer compendia. *Scientific Reports* **6**, 39259 (2016).
- 384 37. Keenan, A. B. *et al.* The Library of Integrated Network-Based Cellular Signatures NIH Program:  
385 System-Level Cataloging of Human Cells Response to Perturbations. *Cell Systems* **6**, 13–24 (2018).
- 386 38. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203  
387 (2018).

- 388 39. Marx, V. Biology: The big challenges of big data. *Nature* **498**, 255–260 (2013).
- 389 40. Abelin, J. G. *et al.* Reduced-representation phosphosignatures measured by quantitative targeted MS  
390 capture cellular states and enable large-scale comparison of drug-induced phenotypes. *Molecular & Cellular*  
391 *Proteomics* mcp.M116.058354 (2016). doi:[10.1074/mcp.M116.058354](https://doi.org/10.1074/mcp.M116.058354)
- 392 41. Bray, M.-A. *et al.* Cell Painting, a high-content image-based assay for morphological profiling using  
393 multiplexed fluorescent dyes. *Nature Protocols* **11**, 1757–1774 (2016).
- 394 42. Bioucas-Dias, J. M. *et al.* Hyperspectral Remote Sensing Data Analysis and Future Challenges. *IEEE*  
395 *Geoscience and Remote Sensing Magazine* **1**, 6–36 (2013).
- 396 43. Berens, P. & Ayhan, M. S. Proprietary data formats block health research. *Nature* **565**, 429 (2019).
- 397 44. Apache Software Foundation. Apache Parquet. <https://parquet.apache.org>
- 398 45. Hipp, D. R. Implementation Limits For SQLite. <https://www.sqlite.org/limits.html>
- 399 46. Furuhashi, S. MessagePack: It’s like JSON. But fast and small. <https://msgpack.org>
- 400 47. Free Software Foundation, Inc. The GNU Awk User’s Guide.  
401 <https://www.gnu.org/software/gawk/manual/gawk.html>
- 402 48. Wickham, H. *Ggplot2: Elegant Graphics for Data Analysis*. (Springer-Verlag New York, 2009).
- 403 49. Wickham, H., Hester, J. & Francois, R. Readr: Read Rectangular Text Data. (2018).
- 404 50. Wickham, H., François, R., Henry, L. & Müller, K. *Dplyr: A Grammar of Data Manipulation*. (2018).
- 405 51. Wilke, C. O. Cowplot: Streamlined Plot Theme and Plot Annotations for ‘ggplot2’. (2017).
- 406 52. McKinney, W. Data Structures for Statistical Computing in Python. in *Proceedings of the 9th Python in*  
407 *Science Conference* 6 (2010).