

1 grepq: A Rust application that quickly filters
2 FASTQ files by matching sequences to a set
3 of regular expressions

4 Nicholas D. Crosbie

5 March 8 2025

6 **grepq: A Rust application that quickly filters**
7 **FASTQ files by matching sequences to a set of**
8 **regular expressions**

9 *Nicholas D. Crosbie, Melbourne Veterinary School, University of Melbourne,*
10 *Parkville, Victoria, Australia*

11 ORCID: 0000-0002-0319-4248

12 March 8, 2025

13 **Keywords:** bioinformatics, FASTQ records, regular expressions,
14 Rust, variants

15 **Abstract**

16 Regular expressions (regex) (Kleene 1951) have been an im-
17 portant tool for finding patterns in biological codes for decades
18 (Hodgman 2000 and citations therein), and unlike fuzzy-finding
19 approaches, do not result in approximate matches. The perfor-
20 mance of regular expressions can be slow, however, especially
21 when searching for matching patterns in large files. *grepq* is
22 a Rust application that quickly filters FASTQ files by matching
23 sequences to a set of regular expressions. *grepq* is designed
24 with a focus on performance and scalability, is easy to install
25 and easy to use, enabling users to quickly filter large FASTQ
26 files, to enumerate named and unnamed variants, to update the
27 order in which patterns are matched against sequences through
28 in-built *tune* and *summarise* commands, and optionally, to output
29 a SQLite file for further sequence analysis. *grepq* is open-source
30 and available on *GitHub* and *Crates.io*.

31 **Statement of need**

32 The ability to quickly filter FASTQ files by matching sequences to
33 a set of regular expressions is an important task in bioinformatics,
34 especially when working with large datasets. The importance and
35 challenge of this task will only grow as sequencing technologies

36 continue to advance and produce ever larger datasets (Katz et
37 al. 2022). The uses cases of *grepq* are diverse, and include pre-
38 processing of FASTQ files before downstream analysis, quality
39 control of sequencing data, and filtering out unwanted sequences.
40 Where decisions need be made quickly, such as in a clinical set-
41 tings (Bachurin et al. 2024), biosecurity (Valdivia-Granda 2012),
42 and wastewater-based epidemiology in support of public health
43 measures (Choi et al. 2018; Sims and Kasprzyk-Hordern 2020;
44 Xylogiannopoulos 2021; Merrett et al. 2024), the ability to quickly
45 filter FASTQ files and enumerate named and unnamed variants by
46 matching sequences to a set of regular expressions is attractive
47 as it circumvents the need for more time-consuming bioinformatic
48 workflows.

49 Regular expressions are a powerful tool for matching sequences,
50 but they can be slow and inefficient when working with large
51 datasets. Furthermore, general purpose tools like *grep* (Free
52 Software Foundation 2023) and *ripgrep* (A. Gallant 2025) are
53 not optimized for the specific task of filtering FASTQ files, and
54 occasionally yield false positives as they scan the entire FASTQ
55 record, including the sequence quality field. Tools such *awk* (Aho,
56 Kernighan, and Weinberger 1988) and *gawk* (Free Software
57 Foundation 2024) can be used to filter FASTQ files without yield-
58 ing false positives, but they are significantly slower than *grepq*

59 and can require the development of more complex scripts to
60 achieve the same result.

61 **Implementation**

62 *grepq* is implemented in Rust, a systems programming language
63 known for its safety features, which help prevent common pro-
64 gramming errors such as null pointer dereferences and buffer over-
65 flows. These features make Rust an ideal choice for implementing
66 a tool like *grepq*, which needs to be fast, efficient, and reliable.

67 Furthermore, *grepq* obtains its performance and reliability, in part,
68 by using the *seq_io* (Schlegel and Seyboldt 2025) and *regex* (Gal-
69 lant et al. 2025b) libraries. The *seq_io* library is a well-tested
70 library for parsing FASTQ files, designed to be fast and efficient,
71 and which includes a module for parallel processing of FASTQ
72 records through multi-threading. The *regex* library is designed
73 to work with regular expressions and sets of regular expressions,
74 and is known to be one of the fastest regular expression libraries
75 currently available (Gallant et al. 2025a). The *regex* library sup-
76 ports Perl-like regular expressions without look-around or backref-
77 erences (documented at https://docs.rs/regex/1.*/regex/#syntax).

78 Further performance gains were obtained by:

- 79 • use of the *RegexSet* struct from the *regex* library to match

80 multiple regular expressions against a sequence in a single
81 pass, rather than matching each regular expression individu-
82 ally (the *RegexSet* is created and compiled once before en-
83 tering any loop that processes the FASTQ records, avoiding
84 the overhead of recompiling the regular expressions for each
85 record)

- 86 • multi-threading to process the records within an input FASTQ
87 file in parallel through use of multiple CPU cores
- 88 • use of the *zlib-ng* backend to the *flate2* library to read and
89 write gzip-compressed FASTQ files, which is faster than the
90 default *miniz_oxide* backend
- 91 • use of an optimised global memory allocator (the *mimalloc*
92 library (Muple, n.d.)) to reduce memory fragmentation and
93 improve memory allocation and deallocation performance
- 94 • buffer reuse to reduce the number of memory allocations and
95 deallocations
- 96 • use of byte slices to avoid the overhead of converting to and
97 from string types
- 98 • in-lining of performance-critical functions
- 99 • use of the *write_all* I/O operation that ensures the data is writ-
100 ten in one go, rather than writing data in smaller chunks

101 **Feature set**

102 *grepq* has the following features:

- 103 • support for presence and absence (inverted) matching of a
104 set of regular expressions
- 105 • IUPAC ambiguity code support (N, R, Y, etc.)
- 106 • support for gzip and zstd compression (reading and writing)
- 107 • JSON support for pattern file input and *tune* and *summarise*
108 command output, allowing named regular expression sets
109 and named regular expressions (pattern files can also be in
110 plain text)
- 111 • the ability to:
 - 112 – set predicates to filter FASTQ records on the header field
113 (= record ID line) using a regular expression, minimum se-
114 quence length, and minimum average quality score (sup-
115 ports Phred+33 and Phred+64)
 - 116 – output matched sequences to one of four formats (includ-
117 ing FASTQ and FASTA)
 - 118 – tune the pattern file and enumerate named and unnamed
119 variants with the *tune* and *summarise* commands: these
120 commands will output a plain text or JSON file with the
121 patterns sorted by their frequency of occurrence in the in-
122 put FASTQ file or gzip-compressed FASTQ file (or a user-
123 specified number of total matches). This can be useful for

124 optimizing the pattern file for performance, for example by
125 removing patterns that are rarely matched and reordering
126 nucleotides within the variable regions of the patterns to
127 improve matching efficiency

- 128 – count and summarise the total number of records and the
129 number of matching records (or records that don't match
130 in the case of inverted matching) in the input FASTQ file
- 131 – bucket matching sequences to separate files named after
132 each regexName with the **–bucket** flag, in any of the four
133 output formats

134 Other than when the **inverted** command is given, output to a
135 SQLite database is supported with the **writeSQL** option. The
136 SQLite database will contain a table called **fastq_data** with the
137 following fields: the fastq record (header, sequence and quality
138 fields), length of the sequence field (length), percent GC content
139 (GC), percent GC content as an integer (GC_int), number of
140 unique tetranucleotides in the sequence (nTN), percent tetranu-
141 cleotide frequency within the sequence (TNF), and a JSON array
142 containing the matched regex patterns, the matches and their po-
143 sition(s) in the FASTQ sequence (variants). If the pattern file was
144 given in JSON format and contained a non-null qualityEncoding
145 field, then the average quality score for the sequence field (av-
146 erage_quality) will also be written. The **–num-tetranucleotides**

147 option can be used to limit the number of tetranucleotides written
148 to the TNF field of the fastq_data SQLite table, these being the
149 most or equal most frequent tetranucleotides in the sequence
150 field of the matched FASTQ records. A summary of the invoked
151 query (pattern and data files) is written to a second table called
152 **query**.

153 Other than when the *tune* or *summarise* command is run, a FASTQ
154 record is deemed to match (and hence provided in the output)
155 when any of the regular expressions in the pattern file match the
156 sequence field of the FASTQ record. Example output of the *tune*
157 command (when given with the **-json-matches** flag) is shown be-
158 low:

```
# For each matched pattern in a search of no more than  
# 20000 matches of a gzip-compressed FASTQ file, print  
# the pattern and the number of matches to a JSON file  
# called matches.json, and include the top three most  
# frequent variants of each pattern, and their respective  
# counts
```

```
grepq --read-gzip 16S-no-iupac.json SRX26365298.fastq.gz \  
tune -n 20000 -c --names --json-matches --variants 3
```

159 Output (abridged) written to matches.json:

```
{  
  "regexSet": {
```



```

    "regex": [
      {
        "regexCount": 2,
        "regexName": "Primer contig 06a",
        "regexString": "[AG]AAT[AT]G[AG]CGGGG",
        "variants": [
          {
            "count": 1,
            "variant": "GAATTGGCGGGG",
            "variantName": "06a-v3"
          },
          {
            "count": 1,
            "variant": "GAATTGACGGGG",
            "variantName": "06a-v1"
          }
        ]
      },
      // matches for other regular expressions...
    ],
    "regexSetName": "conserved 16S rRNA regions"
  }
}

```

160 To output all variants of each pattern, use the `--all` argument, for
161 example:

```
# For each matched pattern in a search of no more than  
# 20000 matches of a gzip-compressed FASTQ file, print  
# the pattern and the number of matches to a JSON file  
# called matches.json, and include all variants of each  
# pattern, and their respective counts. Note that the  
# --variants argument is not given when --all is specified.  
grepq --read-gzip 16S-no-iupac.json SRX26365298.fastq.gz \  
tune -n 20000 -c --names --json-matches --all
```

162 When the count option (`-c`) is given with the *tune* or *summarise*
163 command, *grepq* will count the number of FASTQ records contain-
164 ing a sequence that is matched, for each matching regular expres-
165 sion in the pattern file. If, however, there are multiple occurrences
166 of a given regular expression within a FASTQ record sequence
167 field, *grepq* will count this as one match. To ensure all records are
168 processed, the *summarise* command is used instead of the *tune*
169 command.

170 When the count option (`-c`) is not given as part of the *tune* or *sum-*
171 *marise* command, *grepq* provides the total number of matching
172 FASTQ records for the set of regular expressions in the pattern
173 file.

174 Colorized output for matching regular expressions is not imple-

mented to maximise speed and minimise code complexity, but can be achieved by piping the output to *grep* or *ripgrep* for testing purposes.

Performance

The performance of *grepq* was compared to that of *fqgrep*, *seqkit grep*, *ripgrep*, *grep*, *awk*, and *gawk* using the benchmarking tool *hyperfine*. The test conditions and results are shown in **Table 1**, **Table 2** and **Table 3**.

Table 1: Wall times and speedup of various tools for filtering FASTQ records against a set of regular expressions. Test FASTQ file: SRX26365298.fastq (uncompressed) was 874MB in size, and contained 869,034 records.

tool	wall time (s)		speedup		
	mean	S.D.	× <i>grep</i>	× <i>ripgrep</i>	× <i>awk</i>
<i>grepq</i>	0.192	0.010	1796.76	18.62	863.52
<i>fqgrep</i>	0.338	0.005	1017.61	10.55	489.07
<i>ripgrep</i>	3.568	0.005	96.49	1.00	46.37
<i>seqkit grep</i>	2.885	0.011	119.33	1.24	57.35
<i>grep</i>	344.259	0.545	1.00	0.01	0.48
<i>awk</i>	165.451	1.590	2.08	0.02	1.00
<i>gawk</i>	287.662	1.682	1.20	0.01	0.58

grepq v1.4.0, *fqgrep* v.1.02, *ripgrep* v14.1.1, *seqkit grep* v.2.9.0, *grep* 2.6.0-FreeBSD, *awk* v. 20200816, and *gawk* v.5.3.1. *fqgrep* and *seqkit grep* were run with default settings, *ripgrep* was run with **-B 1 -A 2 --colors 'match:none' --no-line-number**, and *grep* was run with **-B 1 -A 2 --color=never**. *awk* and *gawk* scripts were also configured to output matching records in FASTQ format. The pattern file contained 30 regular expression representing the 12-mers (and

191 their reverse complement) from Table 3 of Martinez-Porchas et al. (2017). The wall times, given in
 192 seconds, are the mean of 10 runs, and S.D. is the standard deviation of the wall times, also given
 193 in seconds.

194 **Table 2:** Wall times and speedup of various tools for filtering gzip-compressed
 195 FASTQ records against a set of regular expressions. Test FASTQ file:
 196 SRX26365298.fastq.gz was 266MB in size, and contained 869,034 records.

tool	wall time (s)		speedup
	mean	S.D.	× ripgrep
<i>grepq</i>	1.703	0.002	2.10
<i>fqgrep</i>	1.834	0.005	1.95
<i>ripgrep</i>	3.584	0.013	1.00

197 Test conditions and tool versions as above, but *grepq* was run with the **–read-gzip** option, *fqgrep*
 198 with the **-Z** option, and *ripgrep* with the **-z** option. SRX26365298.fastq was gzip-compressed using
 199 the *gzip* v.448.0.3 command (Apple Inc. 2019) using default (level 6) settings. The pattern file
 200 contained 30 regular expression representing the 12-mers (and their reverse complement) from
 201 Table 3 of Martinez-Porchas et al. (2017). The wall times, given in seconds, are the mean of 10
 202 runs, and S.D. is the standard deviation of the wall times, also given in seconds.

203 **Table 3:** Wall times and speedup of various tools for filtering FASTQ records
 204 against a set of regular expressions. Test FASTQ file: SRX22685872.fastq was
 205 104GB in size, and contained 139,700,067 records.

206 Test conditions and tool versions as described in the footnote to Table 1. Note that when *grepq* was
 207 run on the gzip-compressed file, a memory resident time for the *grepq* process of 116M as reported
 208 by the *top* command (Apple Inc. 2023c). *fastq-dump* v3.1.1 (Sherry et al. 2012) was used to
 209 download SRX22685872 as a gzip compressed file from the NCBI SRA. The pattern file contained
 210 30 regular expression representing the 12-mers (and their reverse complement) from Table 3 of

tool	wall time (s)		speedup
	mean	S.D.	× ripgrep
<i>Uncompressed</i>			
<i>grepq</i>	26.972	0.244	4.41
<i>fqgrep</i>	50.525	0.501	2.36
<i>ripgrep</i>	119.047	1.227	1.00
<i>gzip-compressed</i>			
<i>grepq</i>	149.172	1.054	0.98
<i>fqgrep</i>	169.537	0.934	0.86
<i>ripgrep</i>	144.333	0.243	1.00

211 Martinez-Porchas et al. (2017). The wall times, given in seconds, are the mean of 10 runs, and
212 S.D. is the standard deviation of the wall times, also given in seconds.

213 Testing

214 The output of *grepq* was compared against the output of *fqgrep*, *seqkit grep*,
215 *ripgrep*, *grep*, *awk* and *gawk*, using the *stat* command (Apple Inc. 2023b),
216 and any difference investigated using the *diff* command (Apple Inc. 2023a).
217 Furthermore, a custom utility, *spikeq* (Crosbie 2024b), was developed to gen-
218 erate synthetic FASTQ files with a known number of records and sequences
219 with user-specified lengths that were spiked with a set of regular expressions a
220 known number of times. This utility was used to test the performance of *grepq*
221 and the aforementioned tools under controlled conditions.

222 Finally, a bash test script (see *examples/test.sh*, available at *grepq*'s Github
223 repository) and a simple Rust CLI application, *predate* (Crosbie 2024a), were
224 developed and utilised to automate system testing, and to monitor for perfor-
225 mance regressions.

226 *grepq* has been tested on macOS 15.0.1 (Apple M1 Max) and Linux Ubuntu
227 20.04.6 LTS (AMD EPYC 7763 64-Core Processor). It may work on other plat-

228 forms, but this has not been tested.

229 **Availability and documentation**

230 *grepq* is open-source and available at *GitHub* (<https://github.com/Rbfinch/grepq>) and *Crates.io* (<https://crates.io/crates/grepq>).

232 Documentation and installation instructions for *grepq* are available at the same
233 GitHub repository, and through the **-h** and **-help** command-line options, which
234 includes a list of all available commands and options, and examples of how to
235 use them. Example pattern files in plain text and JSON format are also provided,
236 as well as test scripts. *grepq* is distributed under the MIT license.

237 **Conclusion**

238 The performance of *grepq* was compared to that of *fqgrep*, *seqkit grep*, *ripgrep*,
239 *grep*, *awk*, and *gawk* using the benchmarking tool *hyperfine*. For an uncom-
240 pressed FASTQ file 874MB in size, containing 869,034 records, *grepq* was
241 significantly faster than the other tools tested, with a speedup of 1797 times
242 relative to *grep*, 864 times relative to *awk*, and 19 times relative to *ripgrep*. For
243 a larger uncompressed FASTQ file (104GB in size, and containing 139,700,067
244 records), *grepq* was 4.4 times faster than *ripgrep* and marginally slower or of
245 equivalent speed to *ripgrep* where the same large file was gzip-compressed.
246 When coupled with its exceptional runtime performance, *grepq*'s feature set
247 make it a powerful and flexible tool for filtering large FASTQ files.

248 Acknowledgements

249 I'm grateful to my family for their patience and support during the development
250 of *grepq*. I would also like to thank the developers of the *seq_io*, *regex*, *mimalloc*
251 and *flate2* libraries for their excellent work, and the developers of the *hyperfine*
252 benchmarking tool for making it easy to compare the performance of different
253 tools. Finally, I would like to thank the authors of the *ripgrep* and *fqgrep* tools
254 for providing inspiration for *grepq*.

255 Conflicts of interest

256 The author declares no conflicts of interest.

257 References

- 258 Aho, Alfred V., Brian W. Kernighan, and Peter J. Weinberger. 1988. *The AWK*
259 *Programming Language*. <https://www.cs.princeton.edu/~bwk/btl.mirror/>.
- 260 Apple Inc. 2019. *The Gzip Command*.
- 261 ———. 2023a. *The Diff Command*.
- 262 ———. 2023b. *The Stat Command*.
- 263 ———. 2023c. *The Top Command*.
- 264 Bachurin, Stanislav S, Mikhail V Yurushkin, Ilya A Slynko, Mikhail E Kletskii,
265 Oleg N Burov, and Dmitriy P Berezovskiy. 2024. "Structural Peculiarities
266 of Tandem Repeats and Their Clinical Significance." *Biochemical and Bio-*
267 *physical Research Communications* 692: 149349.
- 268 Choi, Phil M, Ben J Tschärke, Erica Donner, Jake W O'Brien, Sharon C Grant,
269 Sarit L Kaserzon, Rachel Mackie, et al. 2018. "Wastewater-Based Epidemi-
270 ology Biomarkers: Past, Present and Future." *TrAC Trends in Analytical*
271 *Chemistry* 105: 453–69.

272 Crosbie, Nicholas D. 2024a. “predate: Catch bugs and performance regres-
 273 sions through automated system testing.” <https://github.com/Rbfinch/predate>.
 274 te.

275 ———. 2024b. “spikeq: Generates synthetic FASTQ records free of se-
 276 quences defined by regex patterns, or containing spiked sequences based
 277 on regex patterns.” <https://github.com/Rbfinch/spikeq>.

278 Free Software Foundation. 2023. *GNU Grep 3.11*. Free Software Foundation.
 279 <https://www.gnu.org/software/grep/manual/grep.html>.

280 ———. 2024. *GAWK: Effective AWK Programming: A User’s Guide for GNU*
 281 *Awk, for the 5.3.1*. Free Software Foundation. [https://www.gnu.org/softwa](https://www.gnu.org/software/gawk/manual/gawk.html)
 282 [re/gawk/manual/gawk.html](https://www.gnu.org/software/gawk/manual/gawk.html).

283 Gallant et al. 2025a. “rebar.” <https://github.com/BurntSushi/rebar>.

284 ——— et al. 2025b. “regex.” <https://github.com/rust-lang/regex>.

285 Gallant, Andrew. 2025. “Ripgrep: Recursively Search the Current Directory for
 286 Lines Matching a Pattern.” <https://github.com/BurntSushi/ripgrep>.

287 Hodgman, T. Charles. 2000. “A Historical Perspective on Gene/Protein Func-
 288 tional Assignment.” *Bioinformatics* 16 (1): 10–15.

289 Katz, Kenneth, Oleg Shutov, Richard Lapoint, Michael Kimelman, J Rodney
 290 Brister, and Christopher O’Sullivan. 2022. “The Sequence Read Archive:
 291 A Decade More of Explosive Growth.” *Nucleic Acids Research* 50 (D1):
 292 D387–90.

293 Kleene, SC. 1951. “Representation of Events in Nerve Nets and Finite Au-
 294 tomata.” *CE Shannon and J. McCarthy*.

295 Martinez-Porchas, Marcel, Enrique Villalpando-Canchola, Luis Enrique Ortiz
 296 Suarez, and Francisco Vargas-Albores. 2017. “How Conserved Are the
 297 Conserved 16S-rRNA Regions?” *PeerJ* 5: e3036.

298 Merrett, James E, Monica Nolan, Leon Hartman, Nijoy John, Brianna Flynn,

299 Louise Baker, Christelle Schang, et al. 2024. "Highly Sensitive Wastewater
 300 Surveillance of SARS-CoV-2 Variants by Targeted Next-Generation Ampli-
 301 con Sequencing Provides Early Warning of Incursion in Victoria, Australia."
 302 *Applied and Environmental Microbiology* 90 (8): e01497–23.
 303 Mutiple. n.d. "Mimalloc: A Rust Wrapper over Microsoft's MiMalloc Memory
 304 Allocator."
 305 Schlegel, Markus, and Adrian Seyboldt. 2025. "seq_io: FASTA and FASTQ
 306 parsing and writing in Rust." https://github.com/markschl/seq_io.
 307 Sherry, Stephen, Chunlin Xiao, Kenneth Durbrow, Michael Kimelman, Kurt
 308 Rodarmer, Martin Shumway, and Eugene Yaschenko. 2012. "NCBI Sra
 309 Toolkit Technology for Next Generation Sequence Data." In *Plant and*
 310 *Animal Genome XX Conference (January 14-18, 2012)*. *Plant and Animal*
 311 *Genome*.
 312 Sims, Natalie, and Barbara Kasprzyk-Hordern. 2020. "Future Perspectives
 313 of Wastewater-Based Epidemiology: Monitoring Infectious Disease Spread
 314 and Resistance to the Community Level." *Environment International* 139:
 315 105689.
 316 Valdivia-Granda, Willy A. 2012. "Biodefense Oriented Genomic-Based
 317 Pathogen Classification Systems: Challenges and Opportunities." *Journal*
 318 *of Bioterrorism & Biodefense* 3 (1): 1000113.
 319 Xylogiannopoulos, Konstantinos F. 2021. "Pattern Detection in Multiple
 320 Genome Sequences with Applications: The Case of All SARS-CoV-2
 321 Complete Variants." *bioRxiv*, 2021–04.