

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 29 January 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 January 29, 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 and
25 easy to use, enabling users to quickly filter large FASTQ files, to
26 enumerate named and unnamed variants and update the order
27 in which patterns are matched against sequences through an
28 in-built *tune* command. *grepq* is open-source and available on
29 *GitHub* and *Crates.io*.

30 **Statement of need**

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

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

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

59 achieve the same result.

60 **Implementation**

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

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

77 Further performance gains were obtained by:

- 78 • use of the *RegexSet* struct from the *regex* library to match
79 multiple regular expressions against a sequence in a single

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

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

100 **Feature set**

101 *grepq* has the following features:

- 102 • support for presence and absence (inverted) matching of a
103 set of regular expressions
- 104 • IUPAC ambiguity code support (N, R, Y, etc.)
- 105 • support for gzip and zstd compression (reading and writing)
- 106 • JSON support for pattern file input and *tune* command output,
107 allowing named regular expression sets and named regular
108 expressions (pattern files can also be in plain text)
- 109 • the ability to set predicates to filter FASTQ records on the
110 header field (= record ID line) using a regular expression, min-
111 imum sequence length, and minimum average quality score
112 (supports Phred+33 and Phred+64)
- 113 • the ability to output matched sequences to one of four formats
114 (including FASTQ and FASTA)
- 115 • the ability to tune the pattern file and enumerate named and
116 unnamed variants with the *tune* command: this command
117 will output a plain text or JSON file with the patterns sorted
118 by their frequency of occurrence in the input FASTQ file or
119 gzip-compressed FASTQ file (or a user-specified number of
120 FASTQ records). This can be useful for optimizing the pattern
121 file for performance, for example by removing patterns that
122 are rarely matched and reordering nucleotides within the vari-

123 able regions of the patterns to improve matching efficiency
124 • the ability to count and summarise the total number of records
125 and the number of matching records (or records that don't
126 match in the case of inverted matching) in the input FASTQ
127 file

128 Other than when the *tune* command is run, a FASTQ record is
129 deemed to match (and hence provided in the output) when any
130 of the regular expressions in the pattern file match the sequence
131 field of the FASTQ record. Example output of the *tune* command
132 (when given with the **-json-matches** flag) is shown below:

```
# For each matched pattern in a search of the first  
# 20000 records 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
```

133 Output (abridged) written to matches.json:

```
{  
  "regexSet": {  
    "regex": [  
      {
```

```

    "regexCount": 287,
    "regexName": "Primer contig 06a",
    "regexString": "[AG]AAT[AT]G[AG]CGGGG",
    "variants": [
        {
            "count": 219,
            "variant": "GAATTGACGGGG",
            "variantName": "06a-v1"
        },
        {
            "count": 43,
            "variant": "AAATTGACGGGG",
            "variantName": "06a-v2"
        },
        {
            "count": 21,
            "variant": "GAATTGGCGGGG",
            "variantName": "06a-v3"
        }
    ]
},
// matches for other regular expressions...
],
"regexSetName": "conserved 16S rRNA regions"

```



```
}  
}
```

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

```
# For each matched pattern in a search of the first  
# 20000 records 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
```

136 When the count option (`-c`) is given with the *tune* command, *grepq*
137 will count the number of FASTQ records containing a sequence
138 that is matched, for each matching regular expression in the pat-
139 tern file. If, however, there are multiple occurrences of a given
140 regular expression within a FASTQ record sequence field, *grepq*
141 will count this as one match. When the count option (`-c`) is not
142 given with the *tune* command, *grepq* provides the total number of
143 matching FASTQ records for the set of regular expressions in the
144 pattern file.

145 Colorized output for matching regular expressions is not imple-
146 mented to maximise speed and minimise code complexity, but can

147 be achieved by piping the output to *grep* or *ripgrep* for testing pur-
 148 poses.

149 Performance

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

154 **Table 1:** Wall times and speedup of various tools for filtering FASTQ records
 155 against a set of regular expressions. Test FASTQ file: SRX26365298.fastq
 156 (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

157 *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.
 158 20200816, and *gawk* v.5.3.1. *fqgrep* and *seqkit grep* were run with default settings, *ripgrep* was
 159 run with **-B 1 -A 2 --colors 'match:none' --no-line-number**, and *grep* was run with **-B 1 -A**
 160 **2 --color=never**. *awk* and *gawk* scripts were also configured to output matching records in
 161 FASTQ format. The pattern file contained 30 regular expression representing the 12-mers (and
 162 their reverse compliment) from Table 3 of Martinez-Porchas et al. (2017). The wall times, given in

seconds, are the mean of 10 runs, and S.D. is the standard deviation of the wall times, also given in seconds.

Table 2: Wall times and speedup of various tools for filtering gzip-compressed FASTQ records against a set of regular expressions. Test FASTQ file: 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

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

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

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

177 Test conditions and tool versions as described in the footnote to Table 1. Note that when *grepq* was
178 run on the gzip-compressed file, a memory resident time for the *grepq* process of 116M as reported
179 by the *top* command (Apple Inc. 2023c). *fastq-dump* v3.1.1 (Sherry et al. 2012) was used to
180 download SRX22685872 as a gzip compressed file from the NCBI SRA. The pattern file contained
181 30 regular expression representing the 12-mers (and their reverse compliment) from Table 3 of
182 Martinez-Porchas et al. (2017). The wall times, given in seconds, are the mean of 10 runs, and
183 S.D. is the standard deviation of the wall times, also given in seconds.

184 Testing

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

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

197 *grepq* has been tested on macOS 15.0.1 (Apple M1 Max) and Linux Ubuntu
198 20.04.6 LTS (AMD EPYC 7763 64-Core Processor). It may work on other plat-
199 forms, but this has not been tested.

200 **Availability and documentation**

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

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

208 **Conclusion**

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

219 **Acknowledgements**

220 I'm grateful to my family for their patience and support during the development
221 of *grepq*. I would also like to thank the developers of the *seq_io*, *regex*, *mimalloc*
222 and *flate2* libraries for their excellent work, and the developers of the *hyperfine*

223 benchmarking tool for making it easy to compare the performance of different
224 tools. Finally, I would like to thank the authors of the *ripgrep* and *fgrep* tools
225 for providing inspiration for *grepq*.

226 **Conflicts of interest**

227 The author declares no conflicts of interest.

228 **References**

- 229 Aho, Alfred V., Brian W. Kernighan, and Peter J. Weinberger. 1988. *The AWK*
230 *Programming Language*. <https://www.cs.princeton.edu/~bwk/btl.mirror/>.
- 231 Apple Inc. 2019. *The Gzip Command*.
- 232 ———. 2023a. *The Diff Command*.
- 233 ———. 2023b. *The Stat Command*.
- 234 ———. 2023c. *The Top Command*.
- 235 Bachurin, Stanislav S, Mikhail V Yurushkin, Ilya A Slynko, Mikhail E Kletskii,
236 Oleg N Burov, and Dmitriy P Berezovskiy. 2024. “Structural Peculiarities
237 of Tandem Repeats and Their Clinical Significance.” *Biochemical and Bio-*
238 *physical Research Communications* 692: 149349.
- 239 Choi, Phil M, Ben J Tschärke, Erica Donner, Jake W O’Brien, Sharon C Grant,
240 Sarit L Kaserzon, Rachel Mackie, et al. 2018. “Wastewater-Based Epidemi-
241 ology Biomarkers: Past, Present and Future.” *TrAC Trends in Analytical*
242 *Chemistry* 105: 453–69.
- 243 Crosbie, Nicholas D. 2024a. “predate: Catch bugs and performance regres-
244 sions through automated system testing.” <https://github.com/Rbfinch/predate>.
245 te.
- 246 ———. 2024b. “spikeq: Generates synthetic FASTQ records free of se-

247 quences defined by regex patterns, or containing spiked sequences based
 248 on regex patterns.” <https://github.com/Rbfinch/spikeeq>.
 249 Free Software Foundation. 2023. *GNU Grep 3.11*. Free Software Foundation.
 250 <https://www.gnu.org/software/grep/manual/grep.html>.
 251 ———. 2024. *GAWK: Effective AWK Programming: A User’s Guide for GNU*
 252 *Awk, for the 5.3.1*. Free Software Foundation. <https://www.gnu.org/software/gawk/manual/gawk.html>.
 253
 254 Gallant et al. 2025a. “rebar.” <https://github.com/BurntSushi/rebar>.
 255 ——— et al. 2025b. “regex.” <https://github.com/rust-lang/regex>.
 256 Gallant, Andrew. 2025. “Ripgrep: Recursively Search the Current Directory for
 257 Lines Matching a Pattern.” <https://github.com/BurntSushi/ripgrep>.
 258 Hodgman, T. Charles. 2000. “A Historical Perspective on Gene/Protein Func-
 259 tional Assignment.” *Bioinformatics* 16 (1): 10–15.
 260 Katz, Kenneth, Oleg Shutov, Richard Lapoint, Michael Kimelman, J Rodney
 261 Brister, and Christopher O’Sullivan. 2022. “The Sequence Read Archive:
 262 A Decade More of Explosive Growth.” *Nucleic Acids Research* 50 (D1):
 263 D387–90.
 264 Kleene, SC. 1951. “Representation of Events in Nerve Nets and Finite Au-
 265 tomata.” *CE Shannon and J. McCarthy*.
 266 Martinez-Porchas, Marcel, Enrique Villalpando-Canchola, Luis Enrique Ortiz
 267 Suarez, and Francisco Vargas-Albores. 2017. “How Conserved Are the
 268 Conserved 16S-rRNA Regions?” *PeerJ* 5: e3036.
 269 Merrett, James E, Monica Nolan, Leon Hartman, Nijoy John, Brianna Flynn,
 270 Louise Baker, Christelle Schang, et al. 2024. “Highly Sensitive Wastewater
 271 Surveillance of SARS-CoV-2 Variants by Targeted Next-Generation Ampli-
 272 con Sequencing Provides Early Warning of Incursion in Victoria, Australia.”
 273 *Applied and Environmental Microbiology* 90 (8): e01497–23.

274 Mutiple. n.d. "Mimalloc: A Rust Wrapper over Microsoft's MiMalloc Memory
275 Allocator."

276 Schlegel, Markus, and Adrian Seyboldt. 2025. "seq_io: FASTA and FASTQ
277 parsing and writing in Rust." https://github.com/markschl/seq_io.

278 Sherry, Stephen, Chunlin Xiao, Kenneth Durbrow, Michael Kimelman, Kurt
279 Rodarmer, Martin Shumway, and Eugene Yaschenko. 2012. "NCBI Sra
280 Toolkit Technology for Next Generation Sequence Data." In *Plant and*
281 *Animal Genome XX Conference (January 14-18, 2012)*. *Plant and Animal*
282 *Genome*.

283 Sims, Natalie, and Barbara Kasprzyk-Hordern. 2020. "Future Perspectives
284 of Wastewater-Based Epidemiology: Monitoring Infectious Disease Spread
285 and Resistance to the Community Level." *Environment International* 139:
286 105689.

287 Valdivia-Granda, Willy A. 2012. "Biodefense Oriented Genomic-Based
288 Pathogen Classification Systems: Challenges and Opportunities." *Journal*
289 *of Bioterrorism & Biodefense* 3 (1): 1000113.

290 Xylogiannopoulos, Konstantinos F. 2021. "Pattern Detection in Multiple
291 Genome Sequences with Applications: The Case of All SARS-CoV-2
292 Complete Variants." *bioRxiv*, 2021–04.