

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

4 *Nicholas D. Crosbie, Melbourne Veterinary School, University of Melbourne,*
5 *Parkville, Victoria, Australia*

6 ORCID: 0000-0002-0319-4248

7 January 9, 2025

8 **Keywords:** FASTQ records, regular expressions, Rust, bioinfor-
9 matics

10 **Abstract**

11 Regular expressions (regex) (Kleene 1951) have been an im-
12 portant tool for finding patterns in biological codes for decades
13 (Hodgman 2000 and citations therein), and unlike fuzzy-finding
14 approaches, do not result in approximate matches. The perfor-
15 mance of regular expressions can be slow, however, especially
16 when searching for matching patterns in large files. *grepq* is
17 a Rust application that quickly filters FASTQ files by matching
18 sequences to a set of regular expressions. *grepq* is designed
19 with a focus on performance and scalability, is easy to install
20 and easy to use, enabling users to quickly filter large FASTQ

21 files, and to update the order in which patterns are matched
22 against sequences through an in-built *tune* command. *grepq* is
23 open-source and available on *GitHub* and *Crates.io*.

24 **Statement of need**

25 The ability to quickly filter FASTQ files by matching sequences to
26 a set of regular expressions is an important task in bioinformatics,
27 especially when working with large datasets. The importance and
28 challenge of this task will only grow as sequencing technologies
29 continue to advance and produce ever larger datasets (Katz et
30 al. 2022). The uses cases of *grepq* are diverse, and include pre-
31 processing of FASTQ files before downstream analysis, quality
32 control of sequencing data, and filtering out unwanted sequences.
33 Where decisions need be made quickly, such as in a clinical set-
34 tings (Bachurin et al. 2024), biosecurity (Valdivia-Granda 2012),
35 and wastewater-based epidemiology in support of public health
36 measures (Choi et al. 2018; Sims and Kasprzyk-Hordern 2020;
37 Xylogiannopoulos 2021; Merrett et al. 2024), the ability to quickly
38 filter FASTQ files by matching sequences to a set of regular ex-
39 pressions is attractive as it circumvents the need for more time-
40 consuming bioinformatic workflows.

41 Regular expressions are a powerful tool for matching sequences,
42 but they can be slow and inefficient when working with large

43 datasets. Furthermore, general purpose tools like *grep* (Free
44 Software Foundation 2023) and *ripgrep* (A. Gallant 2025) are
45 not optimized for the specific task of filtering FASTQ files, and
46 occasionally yield false positives as they scan the entire FASTQ
47 record, including the sequence quality field. Tools such *awk* (Aho,
48 Kernighan, and Weinberger 1988) and *gawk* (Free Software
49 Foundation 2024) can be used to filter FASTQ files without yield-
50 ing false positives, but they are significantly slower than *grepq*
51 and can require the development of more complex scripts to
52 achieve the same result.

53 **Implementation**

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

59 Furthermore, *grepq* obtains its performance and reliability, in part,
60 by using the *seq_io* (Schlegel and Seyboldt 2025) and *regex* (Gal-
61 lant et al. 2025b) libraries. The *seq_io* library is a well-tested
62 library for parsing FASTQ files, designed to be fast and efficient,
63 and which includes a module for parallel processing of FASTQ
64 records through multi-threading. The *regex* library is designed

65 to work with regular expressions and sets of regular expressions,
66 and is known to be one of the fastest regular expression libraries
67 currently available (Gallant et al. 2025a). The *regex* library sup-
68 ports Perl-like regular expressions without look-around or backref-
69 erences (documented at https://docs.rs/regex/1.*/regex/#syntax).

70 Further performance gains were obtained by:

- 71 • use of the *RegexSet* struct from the *regex* library to match
72 multiple regular expressions against a sequence in a single
73 pass, rather than matching each regular expression individu-
74 ally (the *RegexSet* is created and compiled once before en-
75 tering any loop that processes the FASTQ records, avoiding
76 the overhead of recompiling the regular expressions for each
77 record)
- 78 • multi-threading to process the records within an input FASTQ
79 file in parallel through use of multiple CPU cores
- 80 • use of the *zlib-ng* backend to the *flate2* library to read and
81 write gzip-compressed FASTQ files, which is faster than the
82 default *miniz_oxide* backend
- 83 • use of an optimised global memory allocator (the *mimalloc*
84 library (Muple, n.d.)) to reduce memory fragmentation and
85 improve memory allocation and deallocation performance
- 86 • buffer reuse to reduce the number of memory allocations and
87 deallocations

- 88 • use of byte slices to avoid the overhead of converting to and
89 from string types
- 90 • in-lining of performance-critical functions
- 91 • use of the *write_all* I/O operation that ensures the data is writ-
92 ten in one go, rather than writing data in smaller chunks

93 **Feature set**

94 *grepq* has the following features:

- 95 • support for presence and absence (inverted) matching of a
96 set of regular expressions
- 97 • IUPAC ambiguity code support (N, R, Y, etc.)
- 98 • gzip support (reading and writing)
- 99 • JSON support for pattern file input and *tune* command output,
100 allowing named regular expression sets and named regular
101 expressions (pattern files can also be in plain text)
- 102 • the ability to set predicates to filter FASTQ records on the
103 header field (= record ID line) using a regular expression, min-
104 imum sequence length, and minimum average quality score
105 (supports Phred+33 and Phred+64)
- 106 • the ability to output matched sequences to one of four formats
107 (including FASTQ and FASTA)
- 108 • the ability to tune the pattern file with the *tune* command: this
109 command will output a plain text or JSON file with the patterns

110 sorted by their frequency of occurrence in the input FASTQ
111 file or gzip-compressed FASTQ file (or a user-specified num-
112 ber of FASTQ records). This can be useful for optimizing the
113 pattern file for performance, for example, by removing pat-
114 terns that are rarely matched

- 115 • the ability to count and summarise the total number of records
116 and the number of matching records (or records that don't
117 match in the case of inverted matching) in the input FASTQ
118 file

119 Other than when the *tune* command is run, a FASTQ record is
120 deemed to match (and hence provided in the output) when any of
121 the regular expressions in the pattern file match the sequence field
122 of the FASTQ record. An example (abridged) output of the *tune*
123 command (when given with the **-json-matches** flag) is shown be-
124 low:

```
{  
  "regexSet": {  
    "regex": [  
      {  
        "regexCount": 287,  
        "regexName": "Primer contig 06a",  
        "regexString": "[AG]AAT[AT]G[AG]CGGGG"  
      },  
    ]  
  }  
}
```

```

{
    "regexCount": 298,
    "regexName": "Primer contig 06aR",
    "regexString": "CCCCG[CT]C[AT]ATT[CT]"
},
{
    "regexCount": 1143,
    "regexName": "Primer contig 03",
    "regexString": "GG[AG][ACGT]GGC[ACGT]GCAG"
}
],
    "regexSetName": "conserved 16S rRNA regions"
}
}

```

125 When the count option (**-c**) is given with the *tune* command, *grepq*
 126 will count the number of FASTQ records containing a sequence
 127 that is matched, for each matching regular expression in the pat-
 128 tern file. If, however, there are multiple occurrences of a given
 129 regular expression within a FASTQ record sequence field, *grepq*
 130 will count this as one match. When the count option (**-c**) is not
 131 given with the *tune* command, *grepq* provides the total number of
 132 matching FASTQ records for the set of regular expressions in the
 133 pattern file.

134 Colorized output for matching regular expressions is not imple-
 135 mented to maximise speed and minimise code complexity, but can
 136 be achieved by piping the output to *grep* or *ripgrep* for testing pur-
 137 poses.

138 Performance

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

143 **Table 1:** Clock times and speedup of various tools for filtering FASTQ records
 144 against a set of regular expressions. Test FASTQ file: SRX26365298.fastq (not
 145 gzip-compressed) was 874MB in size, and contained 869,034 records.

tool	clock time (s)		speedup		
	mean	S.D.	× <i>grep</i>	× <i>ripgrep</i>	× <i>awk</i>
<i>grepq</i>	0.19	0.0021	1814.71	18.74	870.79
<i>fqgrep</i>	0.34	0.01	1010.24	10.43	484.77
<i>ripgrep</i>	3.56	0.01	96.85	1.00	46.48
<i>seqkit grep</i>	122.05	0.90	2.83	0.03	1.36
<i>grep</i>	344.79	1.24	1.00	0.01	0.48
<i>awk</i>	165.45	1.59	2.08	0.02	1.00
<i>gawk</i>	287.66	1.68	1.20	0.01	0.58

146 *grepq* v1.3.6, *fqgrep* v.1.02, *ripgrep* v14.1.1, *seqkit grep* v.2.9.0, *grep* 2.6.0-FreeBSD, *awk* v.
 147 20200816, and *gawk* v.5.3.1. *fqgrep* and *seqkit grep* were run with default settings, *ripgrep* was
 148 run with **—colors 'match:none' —no-line-number**, and *grep* was run with **—color=never**. The
 149 tools were configured to output matching records in FASTQ format. The clock times, given in

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

Table 2: Clock 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	clock time (s)		speedup
	mean	S.D.	× grep
<i>grepq</i>	1.707	0.002	2.09
<i>fqgrep</i>	1.84	0.01	1.94
<i>ripgrep</i>	3.57	0.01	1.00

Test conditions and tool versions as above, but *grepq* was run with the **-x** option, *ripgrep* with the **-z** option, and *grep* with the **-Z** option.

Table 3: Clock times and speedup of various tools for filtering FASTQ records against a set of regular expressions. Test FASTQ file: SRX22685872.fastq (not gzip-compressed) was 104GB in size, and contained 139,700,067 records.

tool	clock time (s)		speedup
	mean	S.D.	× ripgrep
<i>grepq</i>	26.97	0.22	4.38
<i>fqgrep</i>	50.47	0.62	2.34
<i>ripgrep</i>	118.161	1.068	1.00

Test conditions and tool versions as described in the footnote to Table 1. Note that when *grepq* was run under the same conditions but SRX22685872.fastq was gzip-compressed, and a gzip-compressed output was generated, the clock time was 148.01 seconds, with a memory resident time for the *grepq* process of 116M as reported by the *top* command (Apple Inc. 2023c).

164 Testing

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

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

177 *grepq* has been tested on macOS 15.0.1 (Apple M1 Max) and Linux Ubuntu
178 20.04.6 LTS (AMD EPYC 7763 64-Core Processor). It may work on other plat-
179 forms, but this has not been tested.

180 Availability and documentation

181 *grepq* is open-source and available at *GitHub* ([https://github.com/Rbfinch/gre](https://github.com/Rbfinch/grepq)
182 [pq](https://github.com/Rbfinch/grepq)) and *Crates.io* (<https://crates.io/crates/grepq>).

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

188 Conclusion

189 The performance of *grepq* was compared to that of *fqgrep*, *seqkit grep*, *ripgrep*,
190 *grep*, *awk*, and *gawk* using the benchmarking tool *hyperfine*. The results show
191 that *grepq* is significantly faster than the other tools tested, with a speedup of
192 1814.71 times over *grep*, 870.79 times over *awk*, and 18.74 times over *ripgrep*.
193 The performance of *grepq* was also compared to that of *fqgrep* and *ripgrep*
194 when filtering gzip-compressed FASTQ files, with *grepq* being 2.09 times faster
195 than *ripgrep* and 1.94 times faster than *fqgrep*. When coupled with its excep-
196 tional runtime performance, *grepq*'s feature set make it a powerful and flexible
197 tool for filtering large FASTQ files.

198 Acknowledgements

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

205 Conflicts of interest

206 The author declares no conflicts of interest.

207 References

- 208 Aho, Alfred V., Brian W. Kernighan, and Peter J. Weinberger. 1988. *The AWK*
209 *Programming Language*. <https://www.cs.princeton.edu/~bwk/btl.mirror/>.
210 Apple Inc. 2023a. *The Diff Command*.

211 ———. 2023b. *The Stat Command*.

212 ———. 2023c. *The Top Command*.

213 Bachurin, Stanislav S, Mikhail V Yurushkin, Ilya A Slynko, Mikhail E Kletskii,
 214 Oleg N Burov, and Dmitriy P Berezovskiy. 2024. “Structural Peculiarities
 215 of Tandem Repeats and Their Clinical Significance.” *Biochemical and Bio-*
 216 *physical Research Communications* 692: 149349.

217 Choi, Phil M, Ben J Tschärke, Erica Donner, Jake W O’Brien, Sharon C Grant,
 218 Sarit L Kaserzon, Rachel Mackie, et al. 2018. “Wastewater-Based Epidemi-
 219 ology Biomarkers: Past, Present and Future.” *TrAC Trends in Analytical*
 220 *Chemistry* 105: 453–69.

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

224 ———. 2024b. “spikeq: Generates synthetic FASTQ records free of se-
 225 quences defined by regex patterns, or containing spiked sequences based
 226 on regex patterns.” <https://github.com/Rbfinch/spikeq>.

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

229 ———. 2024. *GAWK: Effective AWK Programming: A User’s Guide for GNU*
 230 *Awk, for the 5.3.1*. Free Software Foundation. [https://www.gnu.org/softwa](https://www.gnu.org/software/gawk/manual/gawk.html)
 231 [re/gawk/manual/gawk.html](https://www.gnu.org/software/gawk/manual/gawk.html).

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

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

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

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

238 Katz, Kenneth, Oleg Shutov, Richard Lapoint, Michael Kimelman, J Rodney
 239 Brister, and Christopher O'Sullivan. 2022. "The Sequence Read Archive:
 240 A Decade More of Explosive Growth." *Nucleic Acids Research* 50 (D1):
 241 D387–90.

242 Kleene, SC. 1951. "Representation of Events in Nerve Nets and Finite Au-
 243 tomata." *CE Shannon and J. McCarthy*.

244 Merrett, James E, Monica Nolan, Leon Hartman, Nijoy John, Brianna Flynn,
 245 Louise Baker, Christelle Schang, et al. 2024. "Highly Sensitive Wastewater
 246 Surveillance of SARS-CoV-2 Variants by Targeted Next-Generation Ampli-
 247 con Sequencing Provides Early Warning of Incursion in Victoria, Australia."
 248 *Applied and Environmental Microbiology* 90 (8): e01497–23.

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

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

253 Sims, Natalie, and Barbara Kasprzyk-Hordern. 2020. "Future Perspectives
 254 of Wastewater-Based Epidemiology: Monitoring Infectious Disease Spread
 255 and Resistance to the Community Level." *Environment International* 139:
 256 105689.

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

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