

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

4 *Nicholas D. Crosbie, Melbourne Veterinary School, University of Melbourne,*
5 *Parville, 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 regex patterns. *grepq* is designed with a
19 focus on performance and scalability. *grepq* is easy to install and
20 easy to use, with a simple command-line interface that allows

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

25 **Statement of need**

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

42 Regular expressions are a powerful tool for matching sequences,

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

54 **Implementation**

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

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

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

71 Further performance gains were obtained by:

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

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

94 **Feature set**

95 *grepq* has the following features:

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

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

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

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

```
{  
  "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 regex in the pattern file. If,
 128 however, there are multiple occurrences of a given regex within a
 129 FASTQ record sequence field, *grepq* will count this as one match.
 130 When the count option (**-c**) is not given with the *tune* command,
 131 *grepq* provides the total number of matching FASTQ records for
 132 the set of regex patterns in the pattern file.
 133 Colorized output for matching regex patterns is not implemented

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*, *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: Clock times and speedup of various tools for filtering FASTQ records against a set of regular expressions. Test FASTQ file: SRX26365298.fastq (not 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

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. 20200816, and *gawk* v.5.3.1. *fqgrep* and *seqkit grep* were run with default settings, *ripgrep* was run with **—colors ‘match:none’ —no-line-number**, and *grep* was run with **—color=never**. The 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

150 in seconds.

151 **Table 2:** Clock times and speedup of various tools for filtering gzip-compressed
152 FASTQ records against a set of regular expressions. Test FASTQ file:
153 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

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

156 **Table 3:** Clock times and speedup of various tools for filtering FASTQ records
157 against a set of regular expressions. Test FASTQ file: SRX22685872.fastq (not
158 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

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

163 **Testing**

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

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

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

179 **Availability and documentation**

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

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

187 Discussion and conclusion

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

197 Acknowledgements

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

204 Conflicts of interest

205 The author declares no conflicts of interest.

206 References

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

210 ———. 2023b. *The Stat Command*.

211 ———. 2023c. *The Top Command*.

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

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

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

223 ———. 2024b. “spikeq: Generates synthetic FASTQ records free of se-
 224 quences defined by regex patterns, or containing spiked sequences based
 225 on regex patterns.” <https://github.com/Rbfinch/spikeq>.

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

228 ———. 2024. *GAWK: Effective AWK Programming: A User’s Guide for GNU*
 229 *Awk, for the 5.3.1*. Free Software Foundation. [https://www.gnu.org/softwa](https://www.gnu.org/software/gawk/manual/gawk.html)
 230 [re/gawk/manual/gawk.html](https://www.gnu.org/software/gawk/manual/gawk.html).

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

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

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

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

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

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

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

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

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

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

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

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