

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 22 February 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 February 22, 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 and update the
27 order in which patterns are matched against sequences through
28 in-built *tune* and *summarise* commands. *grepq* is open-source
29 and available on *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* and *summarise*
107 command output, allowing named regular expression sets
108 and named regular expressions (pattern files can also be in
109 plain text)
- 110 • the ability to set predicates to filter FASTQ records on the
111 header field (= record ID line) using a regular expression, min-
112 imum sequence length, and minimum average quality score
113 (supports Phred+33 and Phred+64)
- 114 • the ability to output matched sequences to one of four formats
115 (including FASTQ and FASTA)
- 116 • the ability to tune the pattern file and enumerate named and
117 unnamed variants with the *tune* and *summarise* commands:
118 these commands will output a plain text or JSON file with
119 the patterns sorted by their frequency of occurrence in the
120 input FASTQ file or gzip-compressed FASTQ file (or a user-
121 specified number of total matches). This can be useful for
122 optimizing the pattern file for performance, for example by

123 removing patterns that are rarely matched and reordering nu-
124 cleotides within the variable regions of the patterns to improve
125 matching efficiency

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

133 Other than when the *tune* or *summarise* command is run, a FASTQ
134 record is deemed to match (and hence provided in the output)
135 when any of the regular expressions in the pattern file match the
136 sequence field of the FASTQ record. Example output of the *tune*
137 command (when given with the **--json-matches** flag) is shown be-
138 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
```

139 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"
  }
}

```

140 To output all variants of each pattern, use the `--all` argument, for
 141 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

```

142 When the count option (`-c`) is given with the *tune* or *summarise*
 143 command, *grepq* will count the number of FASTQ records contain-
 144 ing a sequence that is matched, for each matching regular expres-
 145 sion in the pattern file. If, however, there are multiple occurrences
 146 of a given regular expression within a FASTQ record sequence
 147 field, *grepq* will count this as one match. To ensure all records are
 148 processed, the *summarise* command is used instead of the *tune*
 149 command.

150 When the count option (`-c`) is not given as part of the *tune* or *sum-*

151 *marise* command, *grepq* provides the total number of matching
 152 FASTQ records for the set of regular expressions in the pattern
 153 file.

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

158 Performance

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

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

tool	wall time (s)		speedup		
	mean	S.D.	× grep	× ripgrep	× awk
<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

166 *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.
167 20200816, and *gawk* v.5.3.1. *fqgrep* and *seqkit grep* were run with default settings, *ripgrep* was
168 run with **-B 1 -A 2 --colors 'match:none' --no-line-number**, and *grep* was run with **-B 1 -A**
169 **2 --color=never**. *awk* and *gawk* scripts were also configured to output matching records in
170 FASTQ format. The pattern file contained 30 regular expression representing the 12-mers (and
171 their reverse compliment) from Table 3 of Martinez-Porchas et al. (2017). The wall times, given in
172 seconds, are the mean of 10 runs, and S.D. is the standard deviation of the wall times, also given
173 in seconds.

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

tool	wall time (s)		speedup
	mean	S.D.	× <i>ripgrep</i>
<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

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

183 **Table 3:** Wall times and speedup of various tools for filtering FASTQ records
184 against a set of regular expressions. Test FASTQ file: SRX22685872.fastq was
185 104GB in size, and contained 139,700,067 records.

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

Test conditions and tool versions as described in the footnote to Table 1. Note that when *grepq* was run on the gzip-compressed file, a memory resident time for the *grepq* process of 116M as reported by the *top* command (Apple Inc. 2023c). *fastq-dump* v3.1.1 (Sherry et al. 2012) was used to download SRX22685872 as a gzip compressed file from the NCBI SRA. 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.

Testing

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

Finally, a bash test script (see *examples/test.sh*, available at *grepq*'s Github

repository) and a simple Rust CLI application, *predate* (Crosbie 2024a), were developed and utilised to automate system testing, and to monitor for performance regressions.

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

Availability and documentation

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

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

Conclusion

The performance of *grepq* was compared to that of *fqgrep*, *seqkit grep*, *ripgrep*, *grep*, *awk*, and *gawk* using the benchmarking tool *hyperfine*. For an uncompressed FASTQ file 874MB in size, containing 869,034 records, *grepq* was significantly faster than the other tools tested, with a speedup of 1797 times relative to *grep*, 864 times relative to *awk*, and 19 times relative to *ripgrep*. For a larger uncompressed FASTQ file (104GB in size, and containing 139,700,067 records), *grepq* was 4.4 times faster than *ripgrep* and marginally slower or of equivalent speed to *ripgrep* where the same large file was gzip-compressed.

226 When coupled with its exceptional runtime performance, *grepq*'s feature set
227 make it a powerful and flexible tool for filtering large FASTQ files.

228 **Acknowledgements**

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

235 **Conflicts of interest**

236 The author declares no conflicts of interest.

237 **References**

- 238 Aho, Alfred V., Brian W. Kernighan, and Peter J. Weinberger. 1988. *The AWK*
239 *Programming Language*. <https://www.cs.princeton.edu/~bwk/btl.mirror/>.
- 240 Apple Inc. 2019. *The Gzip Command*.
- 241 ———. 2023a. *The Diff Command*.
- 242 ———. 2023b. *The Stat Command*.
- 243 ———. 2023c. *The Top Command*.
- 244 Bachurin, Stanislav S, Mikhail V Yurushkin, Ilya A Slynko, Mikhail E Kletskii,
245 Oleg N Burov, and Dmitriy P Berezovskiy. 2024. "Structural Peculiarities
246 of Tandem Repeats and Their Clinical Significance." *Biochemical and Bio-*
247 *physical Research Communications* 692: 149349.
- 248 Choi, Phil M, Ben J Tscharke, Erica Donner, Jake W O'Brien, Sharon C Grant,

249 Sarit L Kaserzon, Rachel Mackie, et al. 2018. “Wastewater-Based Epidemi-
 250 ology Biomarkers: Past, Present and Future.” *TrAC Trends in Analytical*
 251 *Chemistry* 105: 453–69.

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

255 ———. 2024b. “spikeq: Generates synthetic FASTQ records free of se-
 256 quences defined by regex patterns, or containing spiked sequences based
 257 on regex patterns.” <https://github.com/Rbfinch/spikeq>.

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

260 ———. 2024. *GAWK: Effective AWK Programming: A User’s Guide for GNU*
 261 *Awk, for the 5.3.1*. Free Software Foundation. [https://www.gnu.org/softwa](https://www.gnu.org/software/gawk/manual/gawk.html)
 262 [re/gawk/manual/gawk.html](https://www.gnu.org/software/gawk/manual/gawk.html).

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

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

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

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

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

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

275 Martinez-Porchas, Marcel, Enrique Villalpando-Canchola, Luis Enrique Ortiz

276 Suarez, and Francisco Vargas-Albores. 2017. "How Conserved Are the
 277 Conserved 16S-rRNA Regions?" *PeerJ* 5: e3036.

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

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

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

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

292 Sims, Natalie, and Barbara Kasprzyk-Hordern. 2020. "Future Perspectives
 293 of Wastewater-Based Epidemiology: Monitoring Infectious Disease Spread
 294 and Resistance to the Community Level." *Environment International* 139:
 295 105689.

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

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