

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                   02 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 02, 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

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

48 Regular expressions are a powerful tool for matching sequences,  
49 but they can be slow and inefficient when working with large  
50 datasets. Furthermore, general purpose tools like *grep* (Free  
51 Software Foundation 2023) and *ripgrep* (A. Gallant 2025) are  
52 not optimized for the specific task of filtering FASTQ files, and  
53 occasionally yield false positives as they scan the entire FASTQ  
54 record, including the sequence quality field. Tools such *awk* (Aho,  
55 Kernighan, and Weinberger 1988) and *gawk* (Free Software  
56 Foundation 2024) can be used to filter FASTQ files without yield-  
57 ing false positives, but they are significantly slower than *grepq*  
58 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](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 will  
117 output a plain text or JSON file with the patterns sorted by  
118 their frequency of occurrence in the input FASTQ file or gzip-  
119 compressed FASTQ file (or a user-specified number of total  
120 matches). This can be useful for optimizing the pattern file  
121 for performance, for example by removing patterns that are  
122 rarely matched and reordering nucleotides within the variable

- 123 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 • the ability to bucket matching sequences to separate files  
129 named after each regexName with the `--bucket` flag, in any  
130 of the four output formats

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

```
# 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
```

136 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"
}
}

```



137 To output all variants of each pattern, use the `--all` argument, for  
138 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
```

139 When the count option (`-c`) is given with the *tune* command, *grepq*  
140 will count the number of FASTQ records containing a sequence  
141 that is matched, for each matching regular expression in the pat-  
142 tern file. If, however, there are multiple occurrences of a given  
143 regular expression within a FASTQ record sequence field, *grepq*  
144 will count this as one match. To ensure all records are processed,  
145 the user can supply a large number to the `-n` flag given with the  
146 *tune* command.

147 When the count option (`-c`) is not given with the *tune* command,  
148 *grepq* provides the total number of matching FASTQ records for  
149 the set of regular expressions in the pattern file.

150 Colorized output for matching regular expressions is not imple-  
151 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*, *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* v1.02, *ripgrep* v14.1.1, *seqkit grep* v2.9.0, *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 their reverse complement) 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

182 Test conditions and tool versions as described in the footnote to Table 1. Note that when *grepq* was  
183 run on the gzip-compressed file, a memory resident time for the *grepq* process of 116M as reported  
184 by the *top* command (Apple Inc. 2023c). *fastq-dump* v3.1.1 (Sherry et al. 2012) was used to  
185 download SRX22685872 as a gzip compressed file from the NCBI SRA. The pattern file contained  
186 30 regular expression representing the 12-mers (and their reverse compliment) from Table 3 of  
187 Martinez-Porchas et al. (2017). The wall times, given in seconds, are the mean of 10 runs, and  
188 S.D. is the standard deviation of the wall times, also given in seconds.

## 189 Testing

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

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

202 *grepq* has been tested on macOS 15.0.1 (Apple M1 Max) and Linux Ubuntu  
203 20.04.6 LTS (AMD EPYC 7763 64-Core Processor). It may work on other plat-  
204 forms, but this has not been tested.

## 205 **Availability and documentation**

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

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

## 213 **Conclusion**

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

## 224 **Acknowledgements**

225 I'm grateful to my family for their patience and support during the development  
226 of *grepq*. I would also like to thank the developers of the *seq\_io*, *regex*, *mimalloc*  
227 and *flate2* libraries for their excellent work, and the developers of the *hyperfine*

228 benchmarking tool for making it easy to compare the performance of different  
229 tools. Finally, I would like to thank the authors of the *ripgrep* and *fqgrep* tools  
230 for providing inspiration for *grepq*.

## 231 **Conflicts of interest**

232 The author declares no conflicts of interest.

## 233 **References**

- 234 Aho, Alfred V., Brian W. Kernighan, and Peter J. Weinberger. 1988. *The AWK*  
235 *Programming Language*. <https://www.cs.princeton.edu/~bwk/btl.mirror/>.
- 236 Apple Inc. 2019. *The Gzip Command*.
- 237 ———. 2023a. *The Diff Command*.
- 238 ———. 2023b. *The Stat Command*.
- 239 ———. 2023c. *The Top Command*.
- 240 Bachurin, Stanislav S, Mikhail V Yurushkin, Ilya A Slynko, Mikhail E Kletskii,  
241 Oleg N Burov, and Dmitriy P Berezovskiy. 2024. “Structural Peculiarities  
242 of Tandem Repeats and Their Clinical Significance.” *Biochemical and Bio-*  
243 *physical Research Communications* 692: 149349.
- 244 Choi, Phil M, Ben J Tschärke, Erica Donner, Jake W O’Brien, Sharon C Grant,  
245 Sarit L Kaserzon, Rachel Mackie, et al. 2018. “Wastewater-Based Epidemi-  
246 ology Biomarkers: Past, Present and Future.” *TrAC Trends in Analytical*  
247 *Chemistry* 105: 453–69.
- 248 Crosbie, Nicholas D. 2024a. “predate: Catch bugs and performance regres-  
249 sions through automated system testing.” <https://github.com/Rbfinch/predate>.  
250 te.
- 251 ———. 2024b. “spikeq: Generates synthetic FASTQ records free of se-

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

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

281 Schlegel, Markus, and Adrian Seyboldt. 2025. "seq\_io: FASTA and FASTQ  
282 parsing and writing in Rust." [https://github.com/markschl/seq\\_io](https://github.com/markschl/seq_io).

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

288 Sims, Natalie, and Barbara Kasprzyk-Hordern. 2020. "Future Perspectives  
289 of Wastewater-Based Epidemiology: Monitoring Infectious Disease Spread  
290 and Resistance to the Community Level." *Environment International* 139:  
291 105689.

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

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