# FastqCleaner: an interactive Bioconductor application for quality-control, filtering and trimming of FASTQ files

Leandro Gabriel Roser[*1], Fernán Agüero[1] and Daniel Oscar Sánchez[*1]

## Abstract

**Background**

Exploration and processing of FASTQ files are the first steps in state-of-the-art data analysis workflows of Next Generation Sequencing (NGS) platforms. The large amount of data generated by these technologies has put a challenge in terms of rapid analysis and visualization of sequencing information. Recent integration of the R data analysis platform with web visual frameworks has stimulated the development of user-friendly, powerful, and dynamic NGS data analysis applications.

**Results**

This paper presents *FastqCleaner,* a Bioconductor visual application for both quality-control (QC) and pre-processing of FASTQ files. The interface shows diagnostic information for the input and output data and allows to select a series of filtering and trimming operations in an interactive framework. *FastqCleaner* combines the technology of Bioconductor for NGS data analysis with the data visualization advantages of a web environment.

**Conclusions**

[1] Instituto de Investigaciones Biotecnológicas, Universidad Nacional de San Martín, CONICET, IIB-INTECH, 25 de Mayo y Francia, Buenos Aires, Argentina
* Correspondence: learoser@gmail.com, dsanchez21@gmail.com

20   *FastqCleaner* is an user-friendly, offline-capable tool that enables access to advanced Bioconductor

21   infrastructure. The novel concept of a Bioconductor interactive application that can be used without the

22   need for programming skills, makes *FastqCleaner* a valuable resource for NGS data analysis.

23

24   **Keywords**

25   Bioconductor, FASTQ, Next generation sequencing, R, Shiny, User-friendly tool, Visualization, Web

26   app

27

28   **Background**

29   The advent of Next Generation Sequencing (NGS) technologies has revolutionized genomics,

30   transcriptomics and epigenomics research [1, 2]. The large amount of genetic information produced by

31   these instruments requires suitable data handling and exploration methods. For most common

32   platforms, FASTQ files are the raw starting material for subsequent analyses. A portion of the reads can

33   include adapters or contaminants, the quality of the sequences becomes generally lower towards the

34   end of the reads, and ambiguous base calls may be present. The correction of these and other artifacts

35   are important steps that should be performed before using sequencing reads for mapping or assembly.

36       Bioconductor [3] is a widely used repository based on the R [4] programming language,

37   containing tools devoted to the analysis of high-throughput genomic data. The massive use of these

38   tools is, however, limited by the learning curve that users need to go through to work with customized

39   code routines. Recently, R integration with web tools, in particular JavaScript APIs, has dramatically

40   increased the potential of R to produce more interactive and dynamic experiences of data analysis. This

41   integration is promissory to promote the adoption of R by many researchers for whom learning a

42   programming language has proven to be a prohibitive investment of time and effort.

5                                                           2

43      Here we present *FastqCleaner*, an R package with an offline-capable web application for QC,

44    trimming and filtering of FASTQ files. The tool combines Bioconductor libraries for data analysis and

45    the dynamism of a web application for data visualization.

46

## Implementation

**Application overview**

49    *FastqCleaner* offers the following features:

50    1) Implementation of a local, offline-capable and user-friendly web interface

51    2) Processing of Single-Read (SR) and Paired-End (PE) files

52    3) Dynamic analysis of the input and output files, for customizable sampling size of reads

53    4) Interactive, dynamical exploration and visualization of the data, using cutting-edge technology based

54    on JavaScript and CSS3

55    5) Cross-platform (running in Linux, Mac-OSX and Windows)

56    6) Open source, under GNU GPL (>= 2) license

57

**Program architecture**

59    *FastqCleaner* was developed in R and is distributed as an R package. Data processing is controlled via

60    R functions, that can be also accessed as normal functions from the R console (Additional file 1). These

61    programs make extensive use of the Bioconductor packages *IRanges* [5], *Biostrings* [6] and *ShortRead*

62    [7]. For speed improvement of the routines, C++ code was implemented in R using the *Rcpp* API [8].

63    The web interface included in the package was developed with *Shiny* [9], using JavaScript code written

64    via the *jQuery* API, and CSS3.

65

**Design**

66

67    *FastqCleaner* takes compressed or uncompressed SR or PE files as input (Fig. 1). It accepts files with

68    qualities in both Phred+33 and Phred+64 encoding, detecting Sanger, Solexa and Illumina 1.3+, 1.5+,

69    and > 1.8+ formats. Input files can be processed through a set of independent filters based on either one

70    of the following two principles: 1) *Remotion of a subset of reads that do not meet a given criterion*.

71    This group of filters can remove: a) reads with unknown bases (Ns), b) low complexity sequences, c)

72    duplicated reads, d) reads with length below a threshold quality value, and e) reads with an average

73    quality below a threshold value. 2) *Trimming of individual reads*. This group of filters can trim: a) full

74    and partial adapters, b) 5' regions below a predefined quality threshold, and c) 3' or 5' regions for a

75    fixed nucleotide length. The adapter trimming algorithm extends the methodology of the

76    *trimLRPatterns* function of *Biostrings,* designed to trim on the flanks of reads. For this purpose,

77    *FastqCleaner* includes the *adapter_filter* function, a wrapper of *Biostrings* matching tools. The

78    function is able to trim both adapters present on the flanks or within reads (Fig. 2). Several parameters

79    can be passed to modify the behavior of the tool. These parameters allow, for example, to select a

80    different threshold for the number of mismatches, to take into account the presence of indels, etc.

81          For SR files, *FastqCleaner* sequentially processes a block of reads and writes the resulting post-

82    processed block into the corresponding output file. For PE files, the program uses in each cycle a two-

83    step procedure: first, a block of forward and another of reverse reads are separately processed as in the

84    SR case, and then only those reads present in both post-processed blocks are written into the

85    corresponding output files.

86

**Availability**

87

88    The application and a tutorial are available in Bioconductor at

89    http://bioconductor.org/packages/FastqCleaner/

90

## Installation

92 The application can be installed following the instructions detailed at http://bioconductor.org/packages/

93 FastqCleaner/

94

## Launching the application

96 The application can be launched with the following commands in the R console:

```
97 > library("FastqCleaner")
98 > launch_fqc()
```

99

100 Optionally, when the application is used in RStudio (versions 0.99.878 or higher), a button that allows

101 the direct launch of the application with a single click can be found in the addins menu (Fig. 3).

102

## Results and Discussion

104 The web interface with its three main tabs is described in Fig. 4. The first tab (Fig. 4A) shows the file-

105 selection menu, the available filters, and the run/reset buttons. The selection of files and filters

106 represents the starting point in the *FastqCleaner* workflow. The second tab (Fig. 4B) shows the

107 sequential operations performed on reads after processing. This information consists in the names of

108 the input and output files, and a summary of informative statistics of the reads that passed the filter. The

109 third tab (Figs. 4C, D) shows tables and interactive plots for data diagnostic. Plots can be constructed

110 for both input (original data) and output (post-processed) files. A table with the most frequent k-mers

111 can also be visualized.

112 A comparison of the package with other applications is shown in Table 1. Benchmarking results

113 indicated an excellent performance of *FastqCleaner* in comparison with other pre-processing tools in

8

114     terms of elapsed time. Analysis of SR pre-processing (Fig. 5) showed that these tools can be divided

115     into three groups, in function of significant differences observed in processing speed for routine

116     operations (Tukey HSD test, p < 0.001 for all the three pairwise comparisons). The slowest were

117     *cutadapt* and *FASTX-Toolkit* (group 1), while *AdapterRemoval* and *Trimmomatic* (group 2) were the

118     fastest. *FastqCleaner* showed an intermediate performance, comparable to *Skewer* and *FLEXBAR*

119     (group 3). In PE mode, benchmarking of PE pre-processing operations (Fig. 6) showed that

120     *FastqCleaner* significantly outperforms all other tools for routine operations (Tukey HSD test, *p* <

121     0.001 for pairwise comparisons of *FastqCleaner* vs other individual applications).

122

## Conclusions

124     *FastqCleaner* is a tool with a rich and interactive cutting-edge graphical interface for pre-processing

125     and exploration of SR and PE FASTQ files. Comparison with other available programs in a typical pre-

126     processing scenario of adapter trimming and length filtering, showed an excellent performance of the

127     application for both SR and PE real datasets. The application is made available as an open source

128     license. Coding experience is not required for its use, and is therefore particularly useful for users who

129     are unfamiliar with R programming. Furthermore, all processing happens locally in the user's computer

130     (even if the computer is disconnected from the network), making *FastqCleaner* amenable to run in

131     environments where data confidentiality prevents uploading of files to the cloud.

132         In essence, *FastqCleaner*'s dual capability facilitates both access to the underlying state-of-art

133     Bioconductor infrastructure and to dynamic graphical visualizations in a 100% client-side friendly web

134     environment. This makes *FastqCleaner* a novel technological advance for the analysis of Next

135     Generation Sequencing data.

136

9                                                     6

## Methods

137

138  In order to assess the performance of *FastqCleaner*, we have compared the package with other

139  available pre-processing tools in benchmark tests: *AdapterRemoval* 2.2.2 [10], *cutadapt* 1.14 [11],

140  *FASTX-Toolkit* 0.0.13 [12], *FLEXBAR* 3.0.3 [13], *Skewer* 0.2.2 [14] and *Trimmomatic* 0.36 [15]. The

141  tests (Additional file 2) were conducted for adapter removal and length filtering using SR and PE files,

142  with 22 replicates of each tests for statistical analysis of performance. Processing conditions were

143  standardized by disabling compression of output files and using a single thread. In addition, pre-

144  processing in *FastqCleaner* was performed using a chunk size of 10,000 reads per cycle. For SR

145  processing, we downloaded from SRA the dataset SRR014966, with 14.3 M reads of 36 bp. For PE

146  processing, we downloaded the dataset SRR330569 with 27 M reads of 101 bp. Benchmark tests were

147  conducted in R using a laptop with Linux, a 2.20GHz Intel Core i7 CPU and 16GB of 1600MHz RAM

148  (Additional file 2).

149

## Additional files

150

151  Additional file 1: PDF version of the online tutorial.

152  Additional file 2: R script used in this work for benchmark testing.

153  Additional file 3: Source code of *FastqCleaner* (zip file)

154

## List of abbreviations

155

156  NGS: Next Generation Sequencing

157  SR: Single-Reads

158  PE: Paired-End Reads

10

159

## Availability and requirements

161     **Project name:** FastqCleaner

162     **Project home page:** http://bioconductor.org/packages/FastqCleaner/

163     **Operating system(s):** Platform independent

164     **Programming language:** R, C++, HTML, JavaScript, CSS3

165     **License:** GNU GPL (>= 2)

166

## Declarations

168     **Ethics approval and consent to participate**

169     Not applicable.

170

171     **Consent for publication**

172     Not applicable.

173

174     **Availability of data and materials**

175     *FastqCleaner* is freely available from its Bioconductor home page at http://bioconductor.org/packages/

176     FastqCleaner/ under a GNU GPL (>= 2) license. *FastqCleaner* can be launched on any system that has

177     R installed. An online tutorial is available at the package home page. A PDF version of this tutorial is

178     included as supplemental material (Additional file 1).

179

180     **Competing interests**

181    The authors declare that they have no competing interests.

182

186

187    **Authors' contributions**

188    LGR designed and developed the R package. FA contributed to the improvement of the original design.

189    LGR, FA and DOS wrote the manuscript and tested the package. DOS and FA supervised the project.

190    All authors read and approved the final manuscript.

191

195

196    **Captions to Figures**

197    **Fig. 1** Graphical representation of a typical workflow with *FastqCleaner*, showing the initial selection

198    of FASTQ file(s), processing, and generation of output(s). Diagnostic interactive plots can be

199    constructed for both input and output files. Circular arrows indicate halfway points in the workflow,

200    where different configurations can be selected to re-run the program from there.

201

202    **Fig. 2** Examples for adapter trimming. Pictures show the relative position of an adapter and a read, and

203    the expected result after processing with the *adapter_filter* function of *FastqCleaner*. Dotted lines

204    indicate the portion of the read that will be removed. Arrows show the direction along the read used for

12                                                                        9

205    the program to seek for matches. If one or more matches are found, the function trims the longest

206    subsequence, that contains the matching region plus the rest of the read, in the corresponding trimming

207    direction. *A*: partial adapter on the right + right-trimming of anchored adapter. *B*: partial adapter on the

208    left + left-trimming of anchored adapter. *C*: partial adapter within read + right-trimming. *D,E*: full

209    match between an adapter and a portion of the read + left- (*D*) or right- (*E*) trimming. *F*: multiple

210    matches for a same adapter + left-trimming.

211

212    **Fig. 3** RStudio addins menu, showing the button to launch the *FastqCleaner* application.

213

214    **Fig. 4** Web interface of the *FastqCleaner* application. *A*: first tab, showing an example where a file and

215    a filter are selected. *B*: second tab, showing the processes performed after running the program. *C*: third

216    tab, showing the analysis of the data, in this case for the input FASTQ file. The plot shows the base

217    composition of the sequences. *D*: fourth tab, showing a table with the frequency and the sequence of

218    each duplicated read.

219

220    **Fig. 5** Boxplots for elapsed time (in seconds) for SR adapter trimming and read length filtering.

221

222    **Fig. 6** Boxplots for elapsed time (in seconds) for PE adapter trimming and read length filtering.

223    FASTX-Toolkit is not capable to pre-process PE reads, and hence it is not shown in the plot.

224

## References

226    1. Koboldt D, Steinberg K, Larson D, Wilson R, Mardis E. The next-generation sequencing revolution

227    and its impact on genomics. Cell 2013; 155: 27–38.

13                                                                10

228

229    2. Tripathi R, Sharma P, Chakraborty P, Varadwaj P. Next-generation sequencing revolution through big

230    data analytics. Front Life Sci. 2016; 9: 119–149.

231

232    3. Huber W, Carey V, Gentleman R, Anders S, Carlson M, Carvalho B, Bravo H, Davis S, Gatto L,

233    Girke T, Gottardo R. Orchestrating high-throughput genomic analysis with Bioconductor. Nat Methods.

234    2015; 12: 115–121.

235

236    4. R Core Team. R: A Language and Environment for Statistical Computing. R Foundation for

237    Statistical Computing, Vienna, Austria. 2017.

238

239    5. Lawrence M, Huber W, Pages H, Aboyoun P, Carlson M, Gentleman R, Morgan M, Carey V.

240    Software for Computing and Annotating Genomic Ranges. PLoS Comput Biol. 2013; 9:e1003118.

241

242    6. Pagès H, Aboyoun P, Gentleman R and DebRoy S. Biostrings: String objects representing biological

243    sequences, and matching algorithms. R package version 2.44.2. 2017.

244

245    7. Morgan M, Anders S, Lawrence M, Aboyoun P, Pages H, Gentleman R. ShortRead: a Bioconductor

246    package for input, quality assessment and exploration of high-throughput sequence data.

247    Bioinformatics 2009; 25: 2607–2608.

248

249    8. Eddelbuettel D, Romain F, Allaire J, Chambers J, Bates D, Ushey K. Rcpp: Seamless R and C++

250    integration. J Stat SoftW. 2011; 40: 1-18.

251

252   9. Chang W, Cheng J, Allaire JJ, Xie Y, McPherson J. shiny: Web Application Framework for R. R

253   package version 0.14.2. 2016. https://CRAN.R-project.org/package=shiny

254

255   10. Schubert M, Stinus L, and Ludovic O. AdapterRemoval v2: rapid adapter trimming, identification,

256   and read merging. BMC R Notes. 2016; 9: 88.

257

258   11. Martin, M. (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads.

259   EMBnet journal. 2011; 17:10.

260

261   12.   Hannon   G.   FASTX   Toolkit.   http://hannonlab.cshl.edu/fastx_toolkit/index.html.   Accessed

262   September 5 2017.

263

264   13. Dodt M, Roehr J, Ahmed R, Dieterich C. FLEXBAR—flexible barcode and adapter processing for

265   next-generation sequencing platforms. Biology 2012; 1: 895-905.

266

267   14. Jiang H, Lei R, Ding SW, Zhu S. Skewer: a fast and accurate adapter trimmer for next-generation

268   sequencing paired-end reads. BMC bioinformatics 2014; 15: 182.
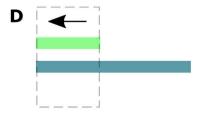
269

270   15. Bolger A, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data.

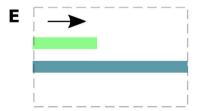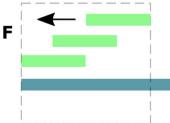271   Bioinformatics 2014; 30: 2114-2120.

```
INPUT
FASTQ
FILE(S)
```

File(s) selection

Filters selection

File processing

Diagnostic plots

```
OUTPUT
FASTQ
FILE(S)
```

File  Edit  Search  View  Macro  Document  Build  Debug  Settings  Tools  Help

Find and Replace
Restore Search
Bookmark Code
In Down to Close Scope?
Insert

**A**

FastqCleaner
A program to clean FASTQ files

Select the operations

1. TRIM BY LENGTH
4. TRIM LOW QUALITY 3' END

2. REMOVE LOW QUALITY DEGREES
5. TRIM 3' BY 5' ADAPTERS

3. REMOVE ADAPTERS
6. TRIM SEQUENCES BY LENGTH

4. REMOVE AVERAGE QUALITY
5. REMOVE LOW SIZE SEQUENCES

Select a file

**B**

FastqCleaner
A program to clean FASTQ files

Input/Output

File input                                     File output

File operations

N filter output:                          Low quality 3' tails filter output:

Low complexity filter output:            Hardlength tails filter output:

Adapter filter output:                    Length filter output:

Average-quality filter output:           Duplicated filter output:

**C**



**D**



Most represented sequences