

GigaScience

SOAPnuke: A MapReduce Acceleration supported Software for integrated Quality Control and Preprocessing of High-Throughput Sequencing Data --Manuscript Draft--

Manuscript Number:	GIGA-D-17-00173	
Full Title:	SOAPnuke: A MapReduce Acceleration supported Software for integrated Quality Control and Preprocessing of High-Throughput Sequencing Data	
Article Type:	Technical Note	
Funding Information:	Science & Technology Bureau of Fujian Province (2013YZ0002-2)	Mr. Lin Fang
	the Joint Project of Natural Science and health Foundation of Fujian Province (2015J01397)	Ms. Chunmei Shi
Abstract:	<p>Quality Control (QC) and preprocessing are essential steps for sequencing data analysis to ensure the accuracy of results. However, existing tools cannot provide a satisfying solution with integrated comprehensive functions, proper architectures and highly-scalable acceleration. In this article, we demonstrate SOAPnuke as an abundant-function tool with 'QC-Preprocess-QC' architecture and MapReduce acceleration framework. Four modules with different preprocessing functions are designed for processing datasets from genomic, small RNA (sRNA), digital gene expression (DGE) and metagenomic experiments respectively. As a workflow-like tool, SOAPnuke centralizes processing functions in one executable and predefine their order to avoid the necessity of reformatting different files when switching tools. Furthermore, the MapReduce framework enables large scalability to distribute all the processing works to the entire computer cluster.</p> <p>We conducted a benchmarking where SOAPnuke and other tools are used to preprocess ~30x NA12878 dataset published by GIAB. The standalone operation of SOAPnuke struck a balance between resource occupancy and performance. When accelerated on 16 working nodes with MapReduce, SOAPnuke achieved 5.7 times of the fastest speed of other tools.</p>	
Corresponding Author:	Lin Fang CHINA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Yuxin Chen	
First Author Secondary Information:		
Order of Authors:	Yuxin Chen	
	Yongsheng Chen	
	Chunmei Shi	
	Zhibo Huang	
	Yong Zhang	
	Shengkang Li	
	Yan Li	
	Jia Ye	
	Chang Yu	

	Zhuo Li
	Xiuqing Zhang
	Jian Wang
	Huanming Yang
	Lin Fang
	Qiang Chen
Order of Authors Secondary Information:	
Opposed Reviewers:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using</p>	Yes

a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist?](#)

1
2
3
4 1 *Technical Note:*

5
6
7 2 **SOAPnuke: A MapReduce Acceleration supported Software for**
8 3 **integrated Quality Control and Preprocessing of**
9 4 **High-Throughput Sequencing Data**

10 5 *Yuxin Chen^{1†}, Yongsheng Chen^{2†}, Chunmei Shi^{3,4,5†}, Zhibo Huang^{1†}, Yong Zhang^{1,6}, Shengkang*
11 6 *Li^{1,6}, Yan Li¹, Jia Ye¹, Chang Yu⁷, Zhuo Li^{8,9}, Xiuqing Zhang¹, Jian Wang^{1,10}, Huanming Yang^{1,10},*
12 7 *Lin Fang^{1,6*} and Qiang Chen^{3,4,5**}*

13
14 8 *1 BGI-Shenzhen, Shenzhen 518083; 2 Geneplus-Beijing, Beijing 102206; 3 Department of Oncology,*
15 9 *Fujian Medical University Union Hospital, Fuzhou 350001; 4 Fujian Key Laboratory of Translational*
16 10 *Cancer Medicine, Fuzhou 350014; 5 Department of Stem Cell Research Institute, Fujian Medical*
17 11 *University Stem Cell Research Institute, Fuzhou 350000; 6 Collaborative Innovation Center of High*
18 12 *Performance Computing, National University of Defense Technology, Changsha 410073; 7 Intel China*
19 13 *Ltd., Shanghai 200336; 8 Guangdong Provincial Hospital of Chinese Medicine, Guangzhou 510120; 9*
20 14 *Department of Surgery, Faculty of Medicine, The Chinese University of Hong Kong, Hong Kong; 10*
21 15 *James D. Watson Institute of Genome Sciences, Hangzhou 310058, China*

22 16 Yuxin C chenyuxin@genomics.cn; Yongsheng C chenys@geneplus.org.cn; CS scmfz@qq.com; ZH
23 17 huangzhibo@genomics.cn; YZ zhangyong2@genomics.cn; SL lishengkang@genomics.cn; YL
24 18 liyan@genomics.cn; JY vejia@genomics.cn; CY chang.yu@intel.com; ZL lzgdphcm@163.com; XZ
25 19 zhangxq@genomics.cn; JW wangjian@genomics.cn; HY yanghm@genomics.cn; LF
26 20 fangl@genomics.cn; QC cqiang8@189.cn

27 21 †Contributed equally

28 22 *First corresponding author

29 23 **Second corresponding author

30 24 **ABSTRACT**

31 25 Quality Control (QC) and preprocessing are essential steps for sequencing data analysis to
32 26 ensure the accuracy of results. However, existing tools cannot provide a satisfying solution with
33 27 integrated comprehensive functions, proper architectures and highly-scalable acceleration. In
34 28 this article, we demonstrate SOAPnuke as an abundant-function tool with ‘QC-Preprocess-QC’
35 29 architecture and MapReduce acceleration framework. Four modules with different
36 30 preprocessing functions are designed for processing datasets from genomic, small RNA (sRNA),
37 31 digital gene expression (DGE) and metagenomic experiments respectively. As a workflow-like

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

32 tool, SOAPnuke centralizes processing functions in one executable and predefine their order to
33 avoid the necessity of reformatting different files when switching tools. Furthermore, the
34 MapReduce framework enables large scalability to distribute all the processing works to the
35 entire computer cluster.

36 We conducted a benchmarking where SOAPnuke and other tools are used to preprocess ~30x
37 NA12878 dataset published by GIAB. The standalone operation of SOAPnuke struck a balance
38 between resource occupancy and performance. When accelerated on 16 working nodes with
39 MapReduce, SOAPnuke achieved 5.7 times of the fastest speed of other tools.

40 **KEYWORDS:** High-throughput sequencing, Quality control, Preprocessing,
41 MapReduce

42 **BACKGROUND**

43 High-throughput sequencing (HTS) instruments have enabled many large-scale studies and
44 generated enormous amount of data [1-3]. However, the presence of low-quality bases, sequence
45 artifacts and sequence contamination can introduce serious negative impact on downstream
46 analyses. Thus, QC and preprocessing of raw data serve as the critical steps to initiate analysis
47 pipelines [4, 5]. QC investigates several statistics of datasets to ensure data quality, and
48 preprocessing trims off undesirable terminal fragments and filters out substandard reads [6].
49 We have conducted a survey on existing 31 tools and widely shared functions are listed in
50 Supplementary Material 1.

51 Existing tools for QC and preprocessing can be divided into two categories according to their
52 structures: toolkit and workflow. Toolkit-like software provides multiple executables such as
53 stats computer, clipper and filtrator [7-15]. In practice, raw data is processed by a few individual

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

54 executables in sequence. Comparatively, workflow-like software offers an integral workflow
55 where functions are performed in predefined order [6, 16-27].

56 However, both categories have their own demerits. When using toolkit-like software, it is
57 complex and error-prone to write additional scripts to wrap executables. Moreover, it consumes
58 much time to generate and read intermediate files, which is hard for acceleration. Besides, the
59 same variables may be computed repetitively. For instance, average quality score of each read is
60 necessary for counting quality score distribution by reads, and filtering reads based on average
61 quality scores. It has to be counted twice if these two functions are implemented by different
62 toolkits.

63 For workflow-like tools, an optimal architecture is required since the orders of functions are
64 fixed. Most of existing tools successively perform QC and preprocessing without complete
65 statistics of preprocessed datasets. If the preprocessing operation is not suitable for a given
66 dataset, the problem will not be revealed until alignment finishes.

67 Datasets sequenced from various samples may require different processing functions or
68 parameters. Existing workflow-like tools mostly support genomics data processing, only a few of
69 them are developed for other types of studies, such as RNA-seq, metagenomics data. For
70 example, RObiNA [22] provides four modules for different RNA sequencing experiments.
71 PrinSeq [6] offers a QC stat, dinucleotide odds ratios, to show how the dataset might be related
72 to other viral/microbial metagenomes. However, there is still no single tool supporting multiple
73 data types.

74 Several tools have made certain progress in overcoming the limitations mentioned above.
75 Galaxy [37] is a web-based platform incorporating various existing toolkit-like software. Users
76 can conveniently concatenate tools into a pipeline on the web interface. NGS QC toolkit [16]

1
2
3
4 77 offers a workflow with QC on un-preprocessed and preprocessed datasets while the
5
6 78 preprocessing functions are too few.

7
8
9
10 79 In terms of software acceleration, only multi-threading is adopted by existing tools [14-16,
11
12 80 24-28]. This approach only works for standalone operation, so it is limited by the maximum
13
14 81 number of processors in one computer server. It may be incompetent when dealing with huge
15
16 82 present and potential volume of sequencing datasets.

17
18
19 83 To solve the problems mentioned above, we developed a workflow-like tool, SOAPnuke, for
20
21 84 integrated QC and preprocessing of large HTS datasets. Similar to NGS QC toolkit, SOAPnuke
22
23 85 performs two-step QC. Frequently used functions of similar tools were integrated in our
24
25
26 86 program, including summary statistics calculation, trimming, filtering and other operations.
27
28 87 Four modules are designed to handle genomic, metagenomic, DGE (Digital Gene Expression
29
30 88 profiling) and sRNA (small RNA) datasets. In addition, SOAPnuke was extended to multiple
31
32 89 working nodes for parallel computing using Hadoop MapReduce framework.

33 34 35 36 90 **METHODS**

37 38 39 91 **QC & PREPROCESSING**

40
41
42
43 92 We developed a two-step procedure to calculate summary statistics of un-preprocessed and
44
45 93 preprocessed data. Basic statistics are comprised of the number of sequences and bases, base
46
47 94 composition, Q20 and Q30, and filtering information. Complex statistics include the
48
49 95 distribution of quality score and base composition distribution for each position. For the former
50
51 96 distribution, Q20 and Q30 for each position were plotted in line chart and the quantiles of the
52
53 97 quality were represented in a boxplot. And for the latter distribution, an overlapping histogram
54
55
56 98 is used to display base composition distribution for each position. The calculations are
57
58 99 conducted by C++ and the figures are generated by R 3.3.2 [38]. An example of the two figures
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

are provided in fig.1. A comprehensive list of statistics available in SOAPnuke is included in Additional File 2. Stats of preprocessed data are compared with some preset thresholds. For per-base quality distribution, a warning will be issued if median score in any position is less than 25, furthermore, a failure will be issued if less than 20. For per-base base composition, a warning will be issued if difference between A and T, or G and C in any position is greater than 10% or a failure will be issued if greater than 20%.

Fig.1. An example of QC complex stats. (a) per-base quality distribution of un-preprocessed paired-end reads. (b) per-base Q20&Q30 of un-preprocessed and preprocessed paired-end reads. (c) per-base base composition distribution of un-preprocessed paired-end reads.

In the step of preprocessing, those undesirable terminal fragments are trimmed off, substandard reads are filtered out, and some transform operations are applied. On both ends of reads, bases of assigned number or of quality lower than threshold will be trimmed off. Sequencing adapters can be aligned, where mismatch is supported while no indel is tolerated, and cut to 3' end. Filtering can be performed on reads with adapter, short length, too many ambiguous bases, low average quality or too many low-quality bases. Those sequencing batches, such as tile for Illumina sequencer and fov for BGI sequencer, with unfavorable sequencing quality can be assigned so the corresponding sequences will be discarded. In addition, reads with identical nucleotides can be deduplicated to keep only one copy. Transformation comprises quality system conversion, interconversion between DNA and RNA, and compression of output with gzip, etc. A comprehensive list of functions mentioned above with their parameters is available in Additional File 3.

MODULES DESIGN

In SOAPnuke, four modules with specified preprocessing functions are set to process different types of data, including General, DGE, sRNA and Meta modules. The previous section shows

1
2
3
4 124 preprocessing functions of General module which can handle most of the DNA re-sequencing
5
6 125 datasets.
7
8
9
10 126 Digital Gene Expression Profiling (DGE) generates single-end read which has a 'CATG' segment
11
12 127 neighboring targeted sequences of 17 base pairs. By default, DGE module will find the targeted
13
14 128 segment and trim off other parts. Moreover, reads with ambiguous bases will be filtered. sRNA
15
16 129 sequences can be contaminated by mRNA during sample preparation. Thus, filtering based on
17
18 130 Poly-A tag is added to mRNA module as polyadenylation is a feature of mRNA data.
19
20
21 131 Metagenomics preprocessing needs less functions than General module. It basically contains
22
23 132 trimming adapters and low-quality bases on both ends, dropping reads with too short length or
24
25 133 too many ambiguous bases. Detailed parameters settings can be accessed in Additional File 3.
26
27

28 134 SOFTWARE FEATURES 29

30
31
32 135 SOAPnuke is written by C++ for good scalability and performance. In terms of cross-platform
33
34 136 compatibility, it can be run on both Linux and Windows.
35
36
37 137 Two paralleled strategies are implemented for acceleration. Multi-threading is developed for
38
39 138 stand-alone operation. Data is cut into blocks of fixed size, and each block is processed by one
40
41 139 thread. This design utilizes multiple cores in a working node. In SOAPnuke, the creation and
42
43 140 allocation of threads are managed by threadpool library, which decreases the overhead of
44
45 141 creating and destroying threads. Hadoop MapReduce is applied to achieve rapid processing in
46
47 142 multi-node cluster. In the mapping phase, each read is kept as a key-value pair, where key is
48
49 143 readID and value is sequence and quality scores. In shuffle phase, the key-value pairs are sorted,
50
51 144 and each pair of paired-end reads is gathered. During the reducing phase, blocks of fixed size are
52
53 145 processed by various threads of multiple nodes, and each block generates an individual result.
54
55 146 After that, it is optional to merge the results into two fastq files.
56
57
58
59
60
61
62
63
64
65

1
2
3
4 147 To prove the effectiveness of the acceleration design, we have conducted a performance tests on
5
6 148 SOAPnuke and other tools. The datasets are published by GIAB [39] and we extracted ~30x for
7
8 149 testing. The access to testing data is available in Addition File 4. In terms of testing environment,
9
10 150 up to 16 nodes were used, each of which has 24 cores of Intel(R) Xeon(R) CPU E5-2620 v4 @
11
12
13 151 2.10GHz and RAM of 128G. To choose operations for testing, we referred to those manuscripts
14
15 152 using SOAPnuke in their works, which are included in Additional File 5. Trimming adapters and
16
17 153 filtering on length and quality were selected for their universality. We chose workflow-like tools
18
19 154 capable of performing these functions, which are Trimmomatic [27], AfterQC [30], BBDuk [31]
20
21 155 and AlignTrimmer [36]. The parameter setting is also available in Addition File 4.
22
23
24

25 156 **RESULTS**

26
27
28 157 In the performance test, we chose three indexes for evaluation: elapsed time, CPU usage and
29
30 158 maximum RAM usage. As shown in table.1, AfterQC is the tool occupying least resources.
31
32
33 159 However, its processing time is too long for practical usage, especially considering we ran the
34
35 160 program with pypy, which is announced to be as 3 times fast as standard python. Among the
36
37 161 remaining tools, SOAPnuke struck an appropriate balance between resources occupancy and
38
39 162 performance. If high performance is demanded, users can choose to run SOAPnuke on multiple
40
41 163 nodes with MapReduce framework. In our testing, 16 nodes can achieve ~32 times acceleration
42
43 164 compared to standalone operation, which is 5.37 times faster than the highest speed of other
44
45 165 tools.
46
47
48
49

	Time (min)	Throughput (read/s)	CPU	Max RAM (GB)
SOAPnuke (1 node 1 thread)	302.7	33947.8	250%	0.62
SOAPnuke (16 nodes)	9.4	1093191.1	640%	50.10

50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Trimmomatic (1 thread)	84.7	121380.1	75%	2.98
Trimmomatic (24 threads)	50.5	203582.1	239%	10.28
BBDuk	57.2	162230.2	259%	11.40
AlienTrimmer	530.2	19076.1	99%	0.54
AfterQC (pypy)	2482.7	4319.1	99%	0.21

Table.1. Evaluation of the data processing performance across SOAPnuke and other 4 tools.

Time, throughput, CPU and maximum memory occupation are presented. For CPU usage, 100% means full load of a single CPU core. Maximum RAM usage means the highest occupancy of RAM during the whole processing.

After the preprocessing, we have conducted the downstream analyses with GATK best practice pipeline (see the description of GATK best practices [<http://www.broadinstitute.org/gatk/guide/best-practices>]). The result shows that the operations of SOAPnuke have similar impact on downstream analysis and details are available in Additional File 4.

DISCUSSION AND CONCLUSION

Quality of data is critical to downstream analysis, which makes it important to use reliable tools for preprocessing. To omit unnecessary IO and computation, workflow-like structure is adopted in SOAPnuke, where QC and preprocessing functions are integrated within an executable.

Compared to most of workflow-like tools, such as PrinSeq [6] and ROBiNA [26], SOAPnuke adds statistics of preprocessed data for better understanding of data. To cope with datasets sequenced from different experiments, four modules are predefined with tailored functions and parameters. In terms of acceleration approach, multi-threading is the sole method adopted by existing tools [14-16, 24-28] but only applicable to single-node operations. SOAPnuke utilizes MapReduce to realize concurrent execution on multi-node operations, where CPU cores of

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

184 multiple nodes can be involved in a single task. It improves the scalability of parallel execution
185 and the applicability to mass data. Our test results indicate that SOAPnuke can achieve ~5.37
186 times faster than the maximum speed of other tools with multi-threading. It is worth
187 mentioning that processing speed is not directly proportional to the number of working nodes,
188 because some procedures like initialization of MapReduce cannot be accelerated as nodes
189 increase, and the burden of communication between nodes aggravates as well.

190 For the future works, we will continue adding functions to various modules. For example, in
191 preprocessing of DGE datasets, filtering out singleton reads is always included [40-42]. For
192 sRNA module, screening out reads based on alignment with some RNA databases (such as tRNA,
193 rRNA and snoRNA) [43,44] is under development. It is also considerable to add statistics like
194 per-read quality distribution and length distribution. In terms of program structure, we are
195 performing refactoring to increase the standalone processing speed.

196 However, we have found two intractable problems with QC and preprocessing respectively
197 among the existing tools. Firstly, in terms of preprocessing, it is difficult to choose optimal
198 parameters for a specific dataset. Datasets from the same experiments and sequencers are tend
199 to share features, so users always select the same parameters for those similar data. The
200 parameters are initially defined based on experiments on a specific dataset or just experience,
201 which already introduces some error. Moreover, even if the parameters are optimal for the
202 tested dataset, they are possibly inappropriate for other data because of other random factors.
203 Thus, the current method is a compromise on processing time. However, it might be a
204 considerable solution that preprocessing settings are automatically adjusted during the
205 processing. Secondly, the QC statistics are of limited help to judge the availability of data. For
206 example, as the threshold of filtering out low-quality reads gets from 0 to 40, all quality-related
207 QC stats will persist rising, but the result of variant calling will be improved then get worse. It is
208 because preprocessing is a procedure required to strike a balance between removing noise and

1
2
3
4 209 keeping useful information, while the balance cannot be reflected by the QC statistics. QC stats
5
6 210 can only indicate whether the filtering condition is too loose instead of too harsh. As a potential
7
8 211 solution, empirical models on similar datasets can be referred to for judging whether the
9
10 212 preprocessing conditions are too harsh. These two problems will also be considered during the
11
12 213 development of following versions.
13
14

15 16 214 **Availability and requirements** 17

18
19 215 Project name: SOAPnuke
20
21

22
23 216 Project home page: <https://github.com/BGI-flexlab/SOAPnuke>
24
25

26 217 Operating system(s): Linux, Windows
27
28

29 218 Programming language: C++
30
31

32
33 219 Requirements: libraries: boost, zlib, log4cplus and openssl; R
34
35

36 220 License: GPL
37
38

39 221 **Abbreviations** 40

41
42
43 222 QC, quality control; HTS, high throughput sequencing; DGE, digital gene expression; sRNA,
44
45 223 small RNA
46
47

48 224 **Declarations** 49

50 51 52 225 **ACKNOWLEDGEMENTS** 53

54
55 226 This research was supported by Collaborative Innovation Center of High Performance
56
57 227 Computing, Critical Patented Project of the Science & Technology Bureau of Fujian Province,
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

228 China (Grant No. 2013YZ0002-2) and the Joint Project of Natural Science and health
229 Foundation of Fujian Province, China (Grant No.2015J01397).

230 **AUTHORS' CONTRIBUTIONS**

231 LF and QC conceived the project. Yuxin C and CS conducted the survey on existing tools for QC
232 and preprocessing. Yuxin C, Yongsheng C, CS, ZH, YZ, SL, JY, ZL, XZ, JW, HY, LF, QC provided
233 feedback on features and functionality. YongSheng C, ZH and SL wrote the standalone version
234 of SOAPnuke. Yuxin C wrote the MapReduce version of SOAPnuke. Yuxin C and ZH performed
235 the above-mentioned test. Yuxin C, YL, CY and LF wrote the manuscript. All authors read and
236 approved the final manuscript.

237 **COMPETING INTERESTS**

238 The authors declare that they have no competing interests.

239 **OPEN ACCESS**

240 This article is distributed under the terms of the Creative Commons Attribution 4.0
241 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits
242 unrestricted use, distribution, and reproduction in any medium, provided you give appropriate
243 credit to the original author(s) and the source, provide a link to the Creative Commons license,
244 and indicate if changes were made. The Creative Commons Public Domain Dedication waiver
245 (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in
246 this article, unless otherwise stated.

247 **ADDITIONAL FILES**

248 Additional file 1: Comparison of features and functions of various tools for QC and
249 preprocessing. (XLSX 41kb)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

250 Additional file 2: Details of QC in SOAPnuke. (PDF 304kb)

251 Additional file 3: Details of preprocessing in SOAPnuke. (PDF 1.6mb)

252 Additional file 4: Details of preprocessing performance test and downstream analyses. (DOCX
253 37kb)

254 Additional file 5: Details of researches involving SOAPnuke. (XLSX 12kb)

255 REFERENCES

- 256 1. Fox S, Filichkin S, Mockler TC. **Applications of ultra-high-throughput sequencing.**
257 Methods Mol Biol. 2009;553:79-108.
- 258 2. Soon WW, Hariharan M, Snyder MP. **High-throughput sequencing for biology and**
259 **medicine.** Mol Syst Biol. 2013;9:640.
- 260 3. Stephens ZD, Lee SY, Faghri F, Campbell RH, Zhai C, Efron MJ, et al. **Big Data:**
261 **Astronomical or Genomical?** PLoS Biol. 2015;13(7):e1002195.
- 262 4. Guo Y, Ye F, Sheng Q, Clark T, Samuels DC. **Three-stage quality control strategies for**
263 **DNA re-sequencing data.** Brief Bioinform. 2014;15(6):879-89.
- 264 5. Zhou X, Rokas A. **Prevention, diagnosis and treatment of high-throughput**
265 **sequencing data pathologies.** Mol Ecol. 2014;23(7):1679-700.
- 266 6. Schmieder R, Edwards R. **Quality control and preprocessing of metagenomic**
267 **datasets.** Bioinformatics. 2011;27(6):863-4.
- 268 7. Moxon S, Schwach F, Dalmay T, Maclean D, Studholme DJ, Moulton V. **A toolkit for**
269 **analysing large-scale plant small RNA datasets.** Bioinformatics. 2008;24(19):2252-3.
- 270 8. Gordon A, Hannon GJ. **Fastx-toolkit. FASTQ/A short-reads preprocessing tools**
271 **(unpublished)** http://hannonlab.cshl.edu/fastx_toolkit. Accessed 1 Mar 2017.

- 1
2
3
4 272 9. Cox MP, Peterson DA, Biggs PJ. **SolexaQA: At-a-glance quality assessment of**
5
6 273 **Illumina second-generation sequencing data.** BMC bioinformatics. 2010 Sep
7
8 274 27;11(1):485.
- 10 275 10. Zhang T, Luo Y, Liu K, Pan L, Zhang B, Yu J, et al. **BIGpre: a quality assessment**
12
13 276 **package for next-generation sequencing data.** Genomics Proteomics Bioinformatics.
14
15 277 2011;9(6):238-44.
- 17 278 11. Aronesty E. **ea-utils: Command-line tools for processing biological sequencing**
18
19 279 **data.** Expression Analysis, Durham, NC. 2011.
- 22 280 12. Yang X, Liu D, Liu F, Wu J, Zou J, Xiao X, et al. **HTQC: a fast quality control toolkit**
23
24 281 **for Illumina sequencing data.** BMC Bioinformatics. 2013;14:33.
- 26 282 13. Li H. **seqtk: Toolkit for processing sequences in FASTA/Q formats.**
27
28 283 <https://github.com/lh3/seqtk>. Accessed 1 Mar 2017.
- 31 284 14. Zhou Q, Su X, Wang A, Xu J, Ning K. **QC-Chain: fast and holistic quality control**
32
33 285 **method for next-generation sequencing data.** PLoS One. 2013;8(4):e60234.
- 35 286 15. Zhou Q, Su X, Jing G, Ning K. **Meta-QC-Chain: comprehensive and fast quality**
36
37 287 **control method for metagenomic data.** Genomics Proteomics Bioinformatics.
38
39 288 2014;12(1):52-6.
- 42 289 16. Patel RK, Jain M. **NGS QC Toolkit: a toolkit for quality control of next generation**
43
44 290 **sequencing data.** PLoS One. 2012;7(2):e30619.
- 46 291 17. Simon A. **FastQC: a quality control tool for high throughput sequence data.**
47
48 292 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/> Accessed 1 Mar 2017
49
- 51 293 18. Schmieder R, Lim YW, Rohwer F, Edwards R. **TagCleaner: Identification and removal**
52
53 294 **of tag sequences from genomic and metagenomic datasets.** BMC Bioinformatics.
54
55 295 2010;11:341.
- 56
57
58
59
60
61
62
63
64
65

- 1
2
3
4 296 19. Falgueras J, Lara AJ, Fernandez-Pozo N, Canton FR, Perez-Trabado G, Claros MG.
5
6 297 **SeqTrim: a high-throughput pipeline for pre-processing any type of sequence**
7
8 298 **read.** BMC Bioinformatics. 2010;11:38.
- 10 299 20. St John, J. **SeqPrep: Tool for stripping adaptors and/or merging paired reads**
11
12 **with overlap into single reads.** <https://github.com/jstjohn/SeqPrep> Accessed 1 Mar
13 300
14 2017
15 301
- 17 302 21. Kong Y. **Btrim: a fast, lightweight adapter and quality trimming program for**
18
19 303 **next-generation sequencing technologies.** Genomics. 2011;98(2):152-3.
- 22 304 22. Lohse M, Bolger AM, Nagel A, Fernie AR, Lunn JE, Stitt M, et al. **RobiNA: a**
23
24 305 **user-friendly, integrated software solution for RNA-Seq-based transcriptomics.**
25
26 306 Nucleic Acids Res. 2012;40(Web Server issue):W622-7.
- 28 307 23. Martin M. **Cutadapt removes adapter sequences from high-throughput**
29
30 308 **sequencing reads.** EMBnet. journal. 2011 May 2;17(1):pp-10.
- 33 309 24. Schubert M, Lindgreen S, Orlando L. **AdapterRemoval v2: rapid adapter trimming,**
34
35 310 **identification, and read merging.** BMC Res Notes. 2016;9:88.
- 37 311 25. Dodt M, Roehr JT, Ahmed R, Dieterich C. **FLEXBAR-Flexible Barcode and Adapter**
38
39 312 **Processing for Next-Generation Sequencing Platforms.** Biology (Basel).
40
41 313 2012;1(3):895-905.
- 44 314 26. Li YL, Weng JC, Hsiao CC, Chou MT, Tseng CW, Hung JH. **PEAT: an intelligent and**
45
46 315 **efficient paired-end sequencing adapter trimming algorithm.** BMC Bioinformatics.
47
48 316 2015;16 Suppl 1:S2.
- 50 317 27. Bolger AM, Lohse M, Usadel B. **Trimmomatic: a flexible trimmer for Illumina**
51
52 318 **sequence data.** Bioinformatics. 2014;30(15):2114-20.
- 55 319 28. Sturm M, Schroeder C, Bauer P. **SeqPurge: highly-sensitive adapter trimming for**
56
57 320 **paired-end NGS data.** BMC Bioinformatics. 2016;17:208.
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

29. Jiang H, Lei R, Ding SW, Zhu S. **Skewer: a fast and accurate adapter trimmer for next-generation sequencing paired-end reads.** BMC Bioinformatics. 2014;15:182.

30. Chen S, Huang T, Zhou Y, Han Y, Xu M, Gu J. **AfterQC: automatic filtering, trimming, error removing and quality control for fastq data.** BMC Bioinformatics. 2017;18(Suppl 3):80.

31. BUSHNELL, Brian. **BBMap: A Fast, Accurate, Splice-Aware Aligner.** Ernest Orlando Lawrence Berkeley National Laboratory, Berkeley, CA (US), 2014.

32. Joshi NA, Fass JN. **Sickle: A sliding-window, adaptive, quality-based trimming tool for FastQ files.** <https://github.com/najoshi/sickle>. Accessed 1 Mar 2017.

33. Perteza, G. **fqtrim: trimming&filtering of next-gen reads.** <https://ccb.jhu.edu/software/fqtrim/>. Access 1 Mar 2017.

34. Vince B. **Scythe: A Bayesian adapter trimmer.** <https://github.com/vsbuffalo/scythe> Access 1 Mar 2017.

35. Leggett RM, Clavijo BJ, Clissold L, Clark MD, Caccamo M. **NextClip: an analysis and read preparation tool for Nextera Long Mate Pair libraries.** Bioinformatics. 2014;30(4):566-8.

36. Criscuolo A, Brisse S. **AlienTrimmer: a tool to quickly and accurately trim off multiple short contaminant sequences from high-throughput sequencing reads.** Genomics. 2013;102(5-6):500-6.

37. Goecks J, Nekrutenko A, Taylor J, Galaxy T. **Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences.** Genome Biol. 2010;11(8):R86.

38. Team RC. **R: A language and environment for statistical computing.** R Foundation for Statistical Computing, Vienna, Austria. 2013.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

39. Zook JM, Catoe D, McDaniel J, Vang L, Spies N, Sidow A, et al. **Extensive sequencing of seven human genomes to characterize benchmark reference materials.** *Sci Data.* 2016;3:160025.

40. Zhang X, Hao L, Meng L, Liu M, Zhao L, Hu F, et al. **Digital gene expression tag profiling analysis of the gene expression patterns regulating the early stage of mouse spermatogenesis.** *PLoS One.* 2013;8(3):e58680.

41. Zhou L, Chen J, Li Z, Li X, Hu X, Huang Y, et al. **Integrated profiling of microRNAs and mRNAs: microRNAs located on Xq27.3 associate with clear cell renal cell carcinoma.** *PLoS One.* 2010;5(12):e15224.

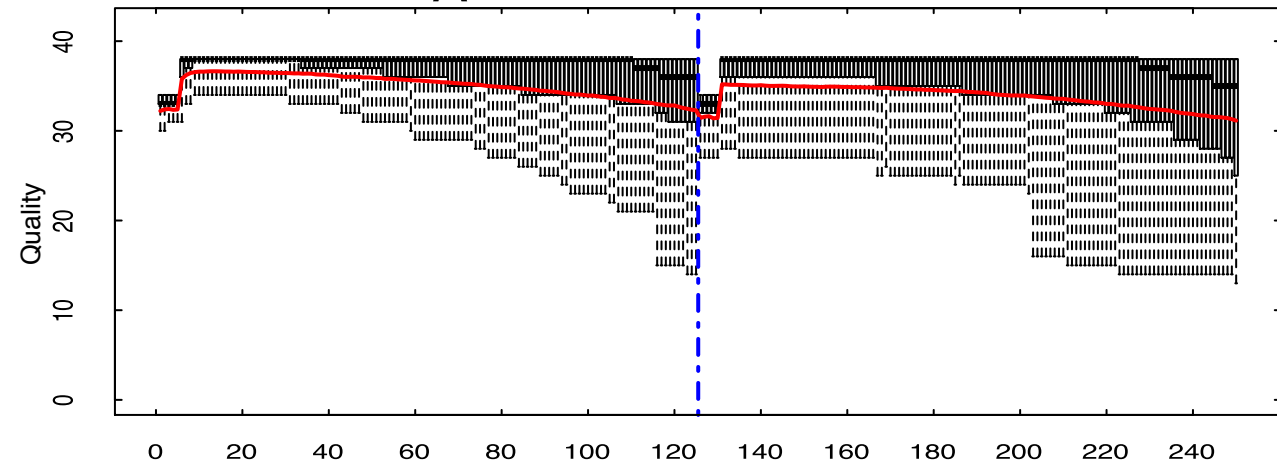
42. Han Y, Zhang X, Wang W, Wang Y, Ming F. **The suppression of WRKY44 by GIGANTEA-miR172 pathway is involved in drought response of Arabidopsis thaliana.** *PLoS One.* 2013;8(11):e73541.

43. Hall AE, Lu WT, Godfrey JD, Antonov AV, Paicu C, Moxon S, et al. **The cytoskeleton adaptor protein ankyrin-1 is upregulated by p53 following DNA damage and alters cell migration.** *Cell Death Dis.* 2016;7:e2184.

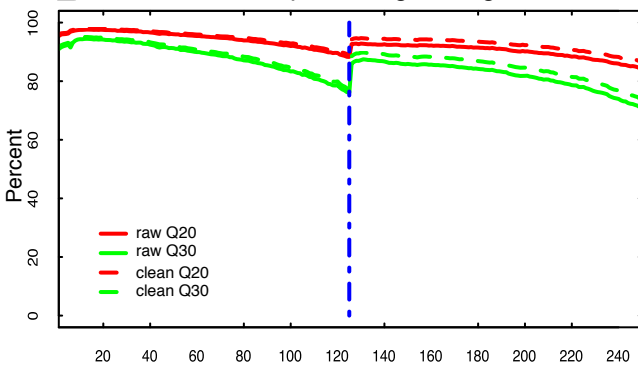
44. Surbanovski N, Brillì M, Moser M, Si-Ammour A. **A highly specific microRNA-mediated mechanism silences LTR retrotransposons of strawberry.** *Plant J.* 2016;85(1):70-82.

Figure.1 [Click here to download Figure Fig.1.pdf](#)

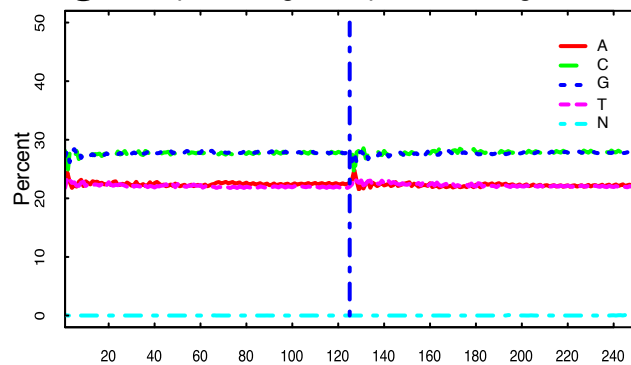
A Base quality along reads



B Q20 Q30 base percentage along reads



C Base percentage composition along reads





Click here to access/download
Supplementary Material
AF1.xlsx





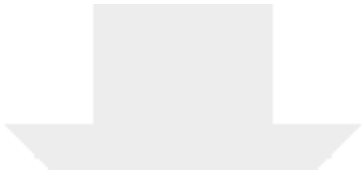
Click here to access/download
Supplementary Material
AF2.pdf






Click here to access/download
Supplementary Material
AF3.pdf





Click here to access/download
Supplementary Material
AF4.docx





Click here to access/download
Supplementary Material
AF5.xlsx

