

SeqyClean

ver. 1.4.7

User Manual

Ilya Y. Zhbannikov

May 31, 2013

Introduction

The main purpose of SeqyClean is to clean next-generation (NG) sequences (“reads”): Roche 454 and Illumina. It provides adapter searching and trimming, quality and poly A/T trimming base on LUCY strategy[1]. SeqyClean accepts both FASTQ and SFF files and also zipped .gz files (zipped FASTQ files only!).

How to download

Follow the link: <https://bitbucket.org/izhbannikov/seqyclean/get/stable.zip>. Save the file under some name.

How to compile

```
cd < path_to_SeqyClean_directory >
make clean
make
```

Usage

Roche 454 sequences

```
./seqyclean [options] < -454 input_filename > < -o output_prefix >
```

Main arguments

< -454 input_filename >	The filename of library to be cleaned. Can be in SFF or FASTQ formats. “-454” tells the program to clean Roche 454 reads
< -o output_prefix >	The files produced will start with the output_prefix followed by some formatted ending (see section “Output files: naming convention”)

Options

-v vector_file	This option is used for vector trimming. If you choose this option, the program assumes the reference genome provided in < vector_file >. This file must be given in FASTA format. Note: vector reference genome(s) must be provided! Example: <code>./seqyclean -v vectors.fa -454 in.sff -o Test</code>
-c file_of_contaminants	This option is used for contaminants screening. If you choose this option, the program assumes the reference genome provided in < file_of_contaminants >. This file must be given in FASTA format. When SeqyClean recognizes contaminants in the sequence, the whole sequence gets discarded. Note: contaminant reference genome(s) must be provided! Example:

	<code>./seqyclean -c contaminants.fa -454 in.sff -o Test</code>
<code>-m file_of_RL_MIDS</code>	This option works in 454 mode only. Use this option to provide your own RLMIDS. SeqyClean will use them and will not use those provided by default. Example: <code>./seqyclean -m file_of_custom_RL_MIDS -454 in.sff -o Test</code>
<code>-k k_mer_size</code>	Use this option in order to specify a size of k-mer. Default size is 15 bases. Example: <code>./seqyclean -k 18 -454 in.sff -o Test</code>
<code>-kc k_mer_size</code>	Special k-mer size for contaminant screening. Use this option only if you want to have different k-mer sizes for contaminant dictionary. Sometimes this option is useful because it prevents false detection of contaminants when program discards too many reads. Example: <code>./seqyclean -kc 25 -454 in.sff -o Test</code>
<code>-f overlap</code>	For Roche 454 only. This option is intended to impose an overlap between two consecutive kmers. By default it is set to 1 bp. Refer to Fig. 1 Example: <code>./seqyclean -f 10 -454 in.sff -o Test</code>
<code>-t number_of_threads</code>	Specifies a number of threads in order to take advantage from using a multicore system. Example: <code>./seqyclean -t 16 -454 in.sff -o Test</code>
<code>-qual max_avg_error max_error_at_ends</code>	LUCY parameters for quality trimming. if "-qual" is set that means you have to provide max_avg_error and max_error_at_ends. Otherwise default values [20 20] will be used. Examples: <code>./seqyclean -qual -454 in.sff -o Test</code> <code>./seqyclean -qual 25 30 -454 in.sff -o Test</code>
<code>--qual_only</code>	Use <code>--qual_only</code> parameter if you want to do only quality trimming. Example: <code>./seqyclean --qual_only -qual -454 in.sff -o Test</code>
<code>--fastq</code>	If input is given in SFF format, by default the output will be also in SFF format. Use this option if you want to have FASTQ format on the output in addition to SFF. Example: <code>./seqyclean --fastq -454 in.sff -o Test</code>
<code>--keep_fastq</code>	Use this option only if you want to keep original FASTQ file from your input SFF Example: <code>./seqyclean --keep_fastq -454 in.sff -o Test</code>
<code>-minimum_read_length <value></code>	Use this option in order to define the minimum number of base pairs when read is still considered as acceptable. If after the cleaning process the read has a length which is less than <code>minimum_read_length</code> parameter, such the read will be discarded. By default, the <code>minimum_read_length</code> is set to 50 base pairs. Example: <code>./seqyclean -minimum_read_length 100 -454 in.sff -o Test</code>
<code>-polyat [cdna] [cerr] [crng]</code>	This option provides trimming of poly A/T tails from nucleotide sequences. Parameters: <code>cdna</code> – tail length (10 by default) <code>cerr</code> – maximum number of errors per tail (3 by default) <code>crng</code> – range to search poly A/T tails (50 by default) Examples: <code>./seqyclean -polyat -454 in.sff -o Test</code> <code>./seqyclean -polyat 12 5 67 poly_test.fastq.gz -o Test.polyAT</code>

Illumina paired- and single-end sequences

Paired-end sequences

```
./seqyclean [options] < -1 input_filename.R1 > < -2 input_filename.R2 > < -o output_prefix >
```

Main arguments

< -1 input_filename.R1 >	The filenames of the library to be cleaned. Must be in FASTQ formats only (the program also accepts zipped (.gz) FASTQ files).
< -2 input_filename.R2 >	
< -o output_prefix >	The files produced will start with the output_prefix followed by some formatted ending (see section "Output files: naming convention")

Single-end sequences

```
./seqyclean [options] < -U input_filename > < -o output_prefix >
```

Main arguments

< -U input_filename >	The filenames of the library to be cleaned. Can be in FASTQ formats only (the program also accepts .gz files).
< -o output_prefix >	The files produced will start with the output_prefix followed by some formatted ending (see section "Output files: naming convention")

Options

-v vector_file	This option is used for vector trimming. If you choose this option, the program assumes the reference genome provided in < vector_file >. This file must be given in FASTA format. Example: ./seqyclean -v vectors.fa -1 R1.fastq.gz -2 R2.fastq.gz -o Test
-c file_of_contaminants	This option is used for contaminants screening. If you choose this option, the program assumes the reference genome provided in < file_of_contaminants >. This file must be given in FASTA format. When SeqyClean recognizes contaminants in the sequence, the whole sequence gets discarded. Example: ./seqyclean -c contaminants.fa -1 R1.fastq.gz -2 R2.fastq.gz -o Test
-k k_mer_size	Use this option in order to specify a size of k-mer. Default size is 15 bases. In Illumina mode this option defines a size of kmer that will be used as a dictionary word size. Example: ./seqyclean -k 14 -1 R1.fastq.gz -2 R2.fastq.gz -o Test
-kc k_mer_size	Special k-mer size for contaminant screening. Use this option only if you want to have different k-mer sizes for contaminant dictionary. Sometimes this option is useful because it prevents false detection of contaminants when program discards too many reads. Example: ./seqyclean -kc 31 -1 R1.fastq.gz -2 R2.fastq.gz -o Test
-qual max_avg_error max_error_at_ends	LUCY parameters for quality trimming. if "-qual" is set that means you have to provide max_avg_error and max_error_at_ends. Otherwise default values [20 20] will be used. Examples: ./seqyclean -qual -1 R1.fastq.gz -2 -1 R2.fastq.gz -o Test ./seqyclean -qual 30 25 -1 R1.fastq.gz -2 R2.fastq.gz -o Test
--qual_only	Use --qual_only parameter if you want to do only quality trimming. Example: ./seqyclean --qual_only -qual -1 R2.fastq.gz -2 R2.fastq.gz -o Test

<code>-minimum_read_length <value></code>	Use this option in order to define the minimum number of base pairs when read is still considered as acceptable. If after cleaning process the read has length which is less than <code>minimum_read_length</code> parameter such read will be discarded. By default, the <code>minimum_read_length</code> is set to 50 base pairs. Note: in this case no adapter/vector/contaminants cleaning is performed. Example: <code>./seqyclean -minimum_read_length 100 -1 R1.fastq.gz -2 R2.fastq.gz -o Test</code>
<code>-polyat [cdna] [cerr] [crng]</code>	This option provides trimming of poly A/T tails from nucleotide sequences. Parameters: <code>cdna</code> – tail length (10 by default) <code>cerr</code> – maximum number of errors per tail (3 by default) <code>crng</code> – range to search poly A/T tails (50 by default) Examples: <code>./seqyclean -polyat -454 in.sff -o Test</code> <code>./seqyclean -polyat 15 4 55 poly_test.fastq.gz -o Test.polyAT</code>

Help

For help please use: `seqyclean -?` or `-help`

Quick examples of usage

Example for 454 reads:

```
./seqyclean -v vectors.fasta -qual 30 25 -454 in.sff -o cleaned_data/Small454Test_cleaned
```

See Figure 2.

Example for Illumina reads:

```
./seqyclean -v vectors.fasta -c contaminants.fasta -qual -1 R1.fastq.gz -2 R2.fastq.gz -o cleaned_data
```

See Figure 3.

Output files: naming conventions

Depending on the given parameters and the cleaning strategy, the name of output file can be different and has the formats described below.

Roche 454

After processing Roche 454 reads, SeqyClean outputs a cleaned file by default in Standard Flowgam Format (SFF) and (if option `--fastq` was chosen) in FASTQ format. Also two report files: `<Prefix>_SummaryStatistics.txt` (which contains information about how many reads were processed, trimmed, discarded and some other information) and `<Prefix>_Report.csv` file which holds the detailed statistics for every read.

Filename
<code><Output_prefix>.sff</code> , <code>.fastq</code> (optionally)
<code><Output_prefix>_Report.tsv</code>
<code><Prefix>_SummaryStatistics.txt</code>
<code><Prefix>_SummaryStatistics.tsv</code>

Illumina

After processing Illumina reads, SeqyClean generates two (shuffled file and file with single-end reads) or three (PE1 and PE2 files that contain paired-end reads and one file with single-end reads) output files in FASTQ format.

Filename
<Output_prefix>_PE1.fastq
<Output_prefix>_PE2.fastq
<Output_prefix>_shuffled.fastq
<Output_prefix>_SE.fastq
<Output_prefix>_PE1.Report.tsv
<Output_prefix>_PE2.Report.tsv
<Prefix>_SummaryStatistics.txt
<Prefix>_SummaryStatistics.tsv

Supported RL MIDs by default

#	Left MID	Right MID	#	Left MID	Right MID
RL1	ACACGACGACT	AGTCGTGGTGT	RL19	ATAGTATACGT	ACGTATAGTAT
RL2	ACACGTAGTAT	ATACTAGGTGT	RL20	CAGTACGTACT	AGTACGTGCTG
RL3	ACACTACTCGT	ACGAGTGGTGT	RL21	CGACGACGCGT	ACGCGTGGTCTG
RL4	ACGACACGTAT	ATACGTGGCGT	RL22	CGACGAGTACT	AGTACTGGTCTG
RL5	ACGAGTAGACT	AGTCTACGCGT	RL23	CGATACTACGT	ACGTAGTGTCTG
RL6	ACGCGTCTAGT	ACTAGAGGCGT	RL24	CGTACGTCGAT	ATCGACGGACG
RL7	ACGTACACACT	AGTGTGTGCGT	RL25	CTACTCGTAGT	ACTACGGGTAG
RL8	ACGTACTGTGT	ACACAGTGCGT	RL26	GTACAGTACGT	ACGTACGGTAC
RL9	ACGTAGATCGT	ACGATCTGCGT	RL27	GTCGTACGTAT	ATACGTAGGAC
RL10	ACTACGTCTCT	AGAGACGGAGT	RL28	GTGTACGACGT	ACGTCTGTGCAC
RL11	ACTATACGAGT	ACTCGTAGAGT	RL29	ACACAGTGAGT	ACTCACGGTGT
RL12	ACTCGCGTCTG	ACGACGGGAGT	RL30	ACACTCATACT	AGTATGGGTGT
RL13	AGACTCGACGT	ACGTCGGGTCT	RL31	ACAGACAGCGT	ACGCTGTGTGT
RL14	AGTACGAGAGT	ACTCTCGGACT	RL32	ACAGACTATAT	ATATAGTGTGT
RL15	AGTACTACTAT	ATAGTAGGACT	RL33	ACAGAGACTCT	AGAGTCTGTGT
RL16	AGTAGACGTCT	AGACGTGCGACT	RL34	ACAGCTCGTGT	ACACGAGGTGT
RL17	AGTCGTACACT	AGTGTAGGACT	RL35	ACAGTGTGCGAT	ATCGACAGTGT
RL18	AGTGTAGTAGT	ACTACTAGACT	RL36	ACGAGCGCGCT	AGCGCGCGCGT

Acknowledgements

This work was supported by IBEST COBRE, grant NIH/NCRR P20RR16448 and the University Research Office at the University of Idaho.

References

- [1] Chou, H. and Holmes, M *DNA sequence cleaning and vector removal* 2001, BMC Bioinformatics, 12, 1093 — 1104.
- [2] <http://www.idtdna.com/pages/products/nextgen/454-adapters>

Contacts

For any questions regarding SeqyClean (i.e. usage, bugs found, performance and so on) please contact Ilya by email: zhba3458@vandals.uidaho.edu. I appreciate every feedback provided by users!
Thank you.

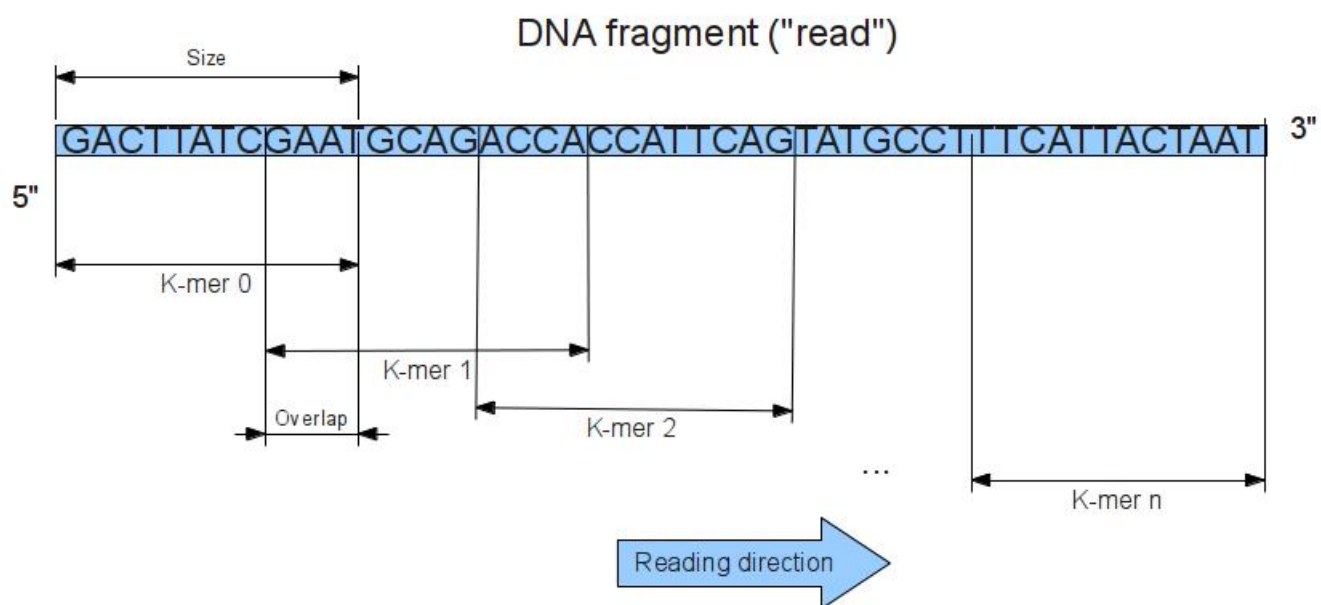


Figure 1: Making a set of consecutive kmers.

```

ilya@kwt:~/bio/app/SeqyClean/bin$ ./seqyclean -qual -v ../test_data/vectors.fasta -454 ../test_data/Small454Test.sff -o ../test/454test
=====Parameters=====
Version: 1.3.12 (2013-03-15)
-----Basic parameters-----
Provided data files :
../test_data/Small454Test.sff
Adapters trimming: YES. Custom RLMIDS: NO. Default RL MIDs will be used.
Vector screening: YES. Vector file provided: ../test_data/vectors.fasta
K-mer size for for vector trimming: 15 bp
Distance between the first bases of two consecutive kmers: 1 bp
Contaminants screening: NO
Quality trimming: YES
Maximum error: 20
Maximum error at ends: 20
-----Output files-----
Output prefix: ../test/454test
Report file: ../test/454test_Report.tsv
Roche output file(s): ../test/454test.sff
-----Other parameters-----
k-mer size: 15 bp
Maximum number of mismatches allowed in alignment: 5
Minimum read length to accept: 50 bp
Distance between the first bases of two consecutive kmers: 1 bp
number_of_threads: 4
=====Starting the process=====
Parsing screening file ../test_data/vectors.fasta
Parsing file: ../test_data/Small454Test.sff...
File is in SFF format, starting conversion...
Conversion finished. Total number of reads read from given file(s): 1000
Running the main pipeline...
Trimming small pieces of vector...
Right ends...
Left ends...
Making clip points...
Making output files...
Making a report...
=====Summary Statistics=====
Reads analyzed: 1000, Bases:988199
Found ->
Left mid tag: 1000, 100%
Right mid tag: 999, 99.9%
# of reads with vector: 212, 21.2%
Reads left trimmed ->
By adapter: 849
By quality: 21
By vector: 128
Average left trim length: 47.1286 bp
Reads right trimmed ->
By adapter: 142
By quality: 804
By vector: 52
Average right trim length: 434.548 bp
Reads discarded: 51 ->
By read length: 51
-----
Reads accepted: 949, %94.9
Average trimmed length: 503.568 bp
=====
Program finished.
Elapsed time = 7.206282e+00 seconds

```

Figure 2: Program output: cleaning Roche 454 reads

```

ilya@kwt:~/bio/app/SeqClean/bin: ./seqclean -qual -v ../test_data/vectors.fasta -1 ../test_data/SmallTestIllumina_R1.fastq.gz -2 ../test_data/SmallTestIllumina_R2.fastq.gz -o ../test/IlluminaTest
=====Parameter=====
Version: 1.3.12 (2013-03-15)
-----Basic parameters-----
Provided data files :
PE1: ../test_data/SmallTestIllumina_R1.fastq.gz, PE2: ../test_data/SmallTestIllumina_R2.fastq.gz
Adapters trimming: YES.
Vector screening: YES. Vector_file provided: ../test_data/vectors.fasta
K-mer_size for for vector trimming: 15
Distance between the first bases of two consecutive kmers: 1
Contaminants screening: NO
Quality trimming: YES
Maximum error: 20
Maximum error at ends: 20
-----Output files-----
Output prefix: ../test/IlluminaTest
Report files: ../test/IlluminaTest_PE1_Report.tsv, ../test/IlluminaTest_PE2_Report.tsv
PE1 file: ../test/IlluminaTest_PE1.fastq
PE2 file: ../test/IlluminaTest_PE2.fastq
Single-end reads: ../test/IlluminaTest_SE.fastq
-----Other parameters-----
Maximum number of mismatches allowed in alignment: 5
Minimum read length to accept: 50
New to old-style Illumina headers: NO
Old-style Illumina: NO
Q-value: 33
=====Starting the process=====
Parsing screening file ../test_data/vectors.fasta
Running the Illumina cleaning process...
Processing files: ../test_data/SmallTestIllumina_R1.fastq.gz, ../test_data/SmallTestIllumina_R2.fastq.gz
i7 adapter in forward first found in the read @MISEQ:6:000000000-A13PY:1:1101:13922:1933 1:N:0:GCCAAT, in the position: 0
i5 adapter in forward first found in the read @MISEQ:6:000000000-A13PY:1:1101:19866:2193 2:N:0:GCCAAT, in the position: 108
=====Summary Statistics=====
PE1 reads analyzed: 31000, Bases:4681000
Found ->
Adapters: 688, 2.21935%
# of reads with vector: 10467, 33.7645%
Reads left trimmed ->
By quality: 21255
By vector: 9745
Average left trim length: 1.38232 bp
Reads right trimmed ->
By adapter: 176
By quality: 4455
By vector: 609
Average right trim length: 13.1112 bp
PE1 reads discarded: 10238
By read length: 10238
-----
PE2 reads analyzed: 31000, Bases:4681000
Found ->
Adapters: 214, 0.690323%
# of reads with vector: 10003, 32.2677%
Reads left trimmed ->
By quality: 21808
By vector: 9186
Average left trim length: 1.52248 bp
Reads right trimmed ->
By quality: 7739
By vector: 662
By adapter: 149
Average right trim length: 16.2237 bp
PE2 reads discarded: 10704
By read length: 10704
-----
-----Summary for PE & SE-----
Pairs kept: 19884, 64.1419%, Bases: 5844269, 62.4254%
Pairs discarded: 9826, 31.6968%, Bases: 2967902, 31.6952%
Single Reads PE1 kept: 878, Bases: 111454
Single Reads PE2 kept: 412, Bases: 47372
Average trimmed length PE1: 147.982 bp
Average trimmed length PE2: 145.936 bp

```

Figure 3: Program output: cleaning paired-end Illumina reads.