

March 4, 2015

SeqyClean User Manual

Ilya Y. Zhbannikov¹, Samuel S. Hunter^{1,2}, Matthew L. Settles^{1,2,3}

¹Institute for Bioinformatics and Evolutionary Studies, University of Idaho, Moscow, ID 83844-3051,

²Department of Biological Sciences, University of Idaho, Moscow, ID 83844-3051,

²Department of Bioinformatics and Computational Biology, University of Idaho, Moscow, ID 83844-3051,

Keywords:

Bioinformatics, Genomics

Correspondence:

Ilya Y. Zhbannikov

Department of Bioinformatics and Computational Biology

University of Idaho

Moscow, ID 83844-30521

Email: zhba3458@vandals.uidaho.edu

Running Title:

SeqyClean User Manual

1 Introduction

We developed SeqyClean – a bioinformatics software pipeline for next-generation sequence cleaning. The first purpose of SeqyClean is to incorporate all aspects of NGS cleaning: adapter, contaminant, poly A/T and quality trimming into a single bioinformatics pipeline. SeqyClean successfully recognizes and removes technological components (adapters, primers, barcodes), contaminants and vector. SeqyClean provides a comprehensive flexible quality trimming by incorporation the LUCY© quality trimming algorithm to remove bad-quality and poly-A/T erroneous data. In addition, SeqyClean has more features: extension paired-end reads by overlap and duplicates removal, which we consider important for genome assembly because it reduces data space by discarding duplicated reads.

2 Installation

2.1 How to download

SeqyClean is an open-source software application available from the Bitbucket for free under this link: <http://bitbucket.org/izhbannikov/seqyclean>. Save the file under some name you wish, unzip and compile:

```
$cd path_to_SeqyClean_directory
$make
```

2.2 Usage

SeqyClean works on SFF files (454, Ion Torrent) and FASTQ Illumina (paired- and single-end reads).

Roche 454 libraries:

```
./seqyclean [options] -454 reads.sff -o output_prefix
```

23 Paired-end Illumina libraries:

24 `./seqyclean [options] -1 R1.fastq -2 R2.fastq -o output_prefix`

26 Single-end libraries:

27 `./seqyclean [options] -U reads.fastq -o output_prefix`

29 **2.3 Options across different technology types**

31 The options that can be used for all library types are shown in Table 1. See library-specific options
32 in the following tables Table 2 (paired-end reads), Table 3 (Roche 454 pyrosequence reads) and
33 Table 4 (single-end FASTQ libraries).

Table 1: Options for all libraries

<code>-v <filename></code>	<p>This option does vector trimming. If you choose this option, the program assumes the file of vector sequences provided in <filename>. This file must be given in FASTA format.</p> <p>Examples:</p> <pre>./seqyclean -v vectors.fa -1 R1.fastq -2 R2.fastq -o Test ./seqyclean -v vectors.fa -U R1.fastq -o Test ./seqyclean -v vectors.fa -454 in.sff -o Test</pre>
<code>-c <filename></code>	<p>This option is used for contaminants screening. If you choose this option, the program assumes the reference genome provided in <filename>. This file must be given in FASTA format. When SeqyClean recognizes contaminants in the sequence, the whole sequence gets discarded. Note: contaminant reference sequences must be provided!</p> <p>Examples:</p> <pre>./seqyclean -v contaminants.fa -1 R1.fastq -2 R2.fastq -o Test ./seqyclean -v contaminants.fa -U R1.fastq -o Test</pre>

	<code>./seqyclean -c contaminants.fa -454 in.sff -o Test</code>
<code>-k <value></code>	Use this option in order to specify a size of k-mer. Default k-mer size is 15 bases.
<code>-kc <value></code>	Special k-mer size for contaminant screening. Use this option only if you want to have different k-mer sizes for contaminant dictionary.
<code>-qual [mae mee -w0 <value> -w1 <value>]</code>	Quality trimming. Default values for mae (maximum average error) and mee (maximum error at ends) are [0.01 0.01]. "w0" and "w1" are window parameters. Examples: <code>./seqyclean -1 R1.fastq -2 R2.fastq -o Test -qual</code> <code>./seqyclean -qual 0.012 -w0 40 -w1 5 -U R1.fastq -o Test</code> <code>./seqyclean -qual 0.025 0.0030 -454 in.sff -o Test</code>
<code>-bracket [bracket length] [max avg error]</code>	Bracket parameters: minimum length (default=10) and maximum average error (default=0.794 or 1 phred) - these maximum average error values means that checking for bracket error is OFF)
<code>-window window_size max_avg_error [window_size max_avg_error ...]</code>	Parameters for window trimming. By default two windows are used: large window, 50 bp long, with maximum average error of 0.794 and small window, 10 bp long, with maximum average error of 0.794. By default checking for error at this stage of quality trimming algorithm is OFF.
<code>-minlen value</code>	Use this option <code>-minlen</code> 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>-minlen</code> parameter, the read will be discarded. By default, the <code>-minlen</code> is set to 50 base pairs. Example: <code>./seqyclean -minlen 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. cdna – tail length (10 by default); cerr – maximum number of errors per tail (3 by default); crng – range to search poly A/T tails (50 by default) Examples: <code>./seqyclean -polyat -1 R1.fastq -2 R2.fastq -o Test</code> <code>./seqyclean -polyat 12 5 120 -U R1.fastq -o Test</code>

	<code>./seqyclean -polyat -454 in.sff -o Test</code>
<code>-dup [startdw] [sizedw] [maxdup]</code>	<p>This option provides duplicates screening.</p> <p><code>startdw</code> – search starting position (10 by default); <code>sizedw</code> – size of window (35 by default); <code>maxdup</code> – maximum number of duplicates (3 by default)</p> <p>Examples:</p> <pre>./seqyclean -dup -1 R1.fastq -2 R2.fastq -o Test ./seqyclean -dup -sizedw 50 -U R1.fastq -o Test ./seqyclean -dup -startdw 5 -sizedw 30 -maxdup 12 -454 in.sff -o Test</pre>
<code>-verbose</code>	Verbose output, default=off.
<code>-detrep</code>	Generate detailed report for each read, default=off.
<code>-no_adapter_trim</code>	This option turns off adapter trimming. Default=off.

Table 2: Illumina paired-end libraries

<code>-shuffle</code>	<p>With this option SeqyClean will combine output paired-end libraries into one single file named <code><output_prefix>_shuffled.fastq</code>. However, SeqyClean still does keep single-end reads (reads without corresponding pairs) in <code><output_prefix>_SE.fastq</code> file.</p> <p>Example:</p> <pre>./seqyclean -shuffle -1 R1.fastq -2 R2.fastq -o Test</pre>
<code>-at <value></code>	This option sets the similarity threshold for adapter trimming by overlap. By default its value is set to 0.75.
<code>-alen <value></code>	This option sets the maximum adapter length for adapter trimming by overlap. By default its value is set to 60 bases.
<code>-overlap <value></code>	This option turns on merging overlapping paired-end reads and <code><value></code> is the minimum overlap length. By default the minimum overlap length is 16 base pairs.
<code>-i64</code>	Turns on 64-quality base, default = off.
<code>-new2old</code>	A switch to fix read IDs, default=off (As is detailed in: http://contig.wordpress.com/2011/09/01/newbler-input-iii-a-quick-fix-for-the-new-illumina-fastq-header/)

2.4 Description of seqyclean output

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

2.4.1 SFF (454, Ion Torrent)

- Output_prefix.sff , .fastq (optionally)
- Output_prefix_Report.tsv - if `-det rep` flag is on.
- Prefix_SummaryStatistics.txt
- Prefix_SummaryStatistics.tsv

Table 3: Roche 454 pyrosequence libraries

<code>-t <value></code>	Number of threads (not yet applicable to Illumina mode), default=4.
<code>-fastq</code>	Output in FASTQ format, default=off.
<code>-fasta</code>	Output in FASTA format, default=off.
<code>-m <filename></code>	Using custom barcodes, default=off. <code>filename</code> - a path to a FASTA-file with custom barcodes.
<code>-d <value></code>	This option <code>-d</code> is intended to tweak an overlap between two consecutive k-mers. By default the length of overlap it is set to 1 bp. Example: <code>./seqyclean -d 10 -454 in.sff -o Test</code>

Table 4: Single-end FASTQ libraries

<code>-U <filename></code>	Turns on single-end mode.
<code>-i64</code>	Turns on 64-quality base, default = off.
<code>-new2old</code>	A switch to fix read IDs, default=off (As is detailed in: http://contig.wordpress.com/2011/09/01/newbler-input-iii-a-quick-fix-for-the-new-illumina-fastq-header/)

2.5 FASTQ

After processing FASTQ reads, SeqyClean generates PE1 and PE2 files that contain paired-end reads, SE file with single-end reads OR 'shuffled' file and file with single-end reads (SE) if `-shuffle` flag was set. output files in FASTQ format.

- `Output_prefix_PE1.fastq`
- `Output_prefix_PE2.fastq`
- `Output_prefix_shuffled.fastq` (if `-shuffle` flag was set)
- `Output_prefix_SE.fastq`
- `Output_prefix_PE1_Report.tsv` (if `-det rep` flag was set)
- `Output_prefix_PE2_Report.tsv` (if `-det rep` flag is on)
- `Prefix_SummaryStatistics.txt`
- `Prefix_SummaryStatistics.tsv`

2.6 Workflow

The general workflow diagram of SeqyClean is shown in Figure 1 and described below. The workflow consists of several atomic steps: (1) Input data pre-processing; (2) Trimming poly A/T tails; (3) Vector and contaminants trimming; (4) Adapter trimming; (5) Quality trimming; (6) Extension by overlap; (7) PCR duplicates removal; (8) Establishing clip points; (9) Generating output files and summary statistics. Stages 2, 3, 4, 5, 6, 7 are optional depending on chosen cleaning strategy.

2.6.1 Supported RLMIDs

The set of supported Roche 454 RL MIDs is shown in Table 5.

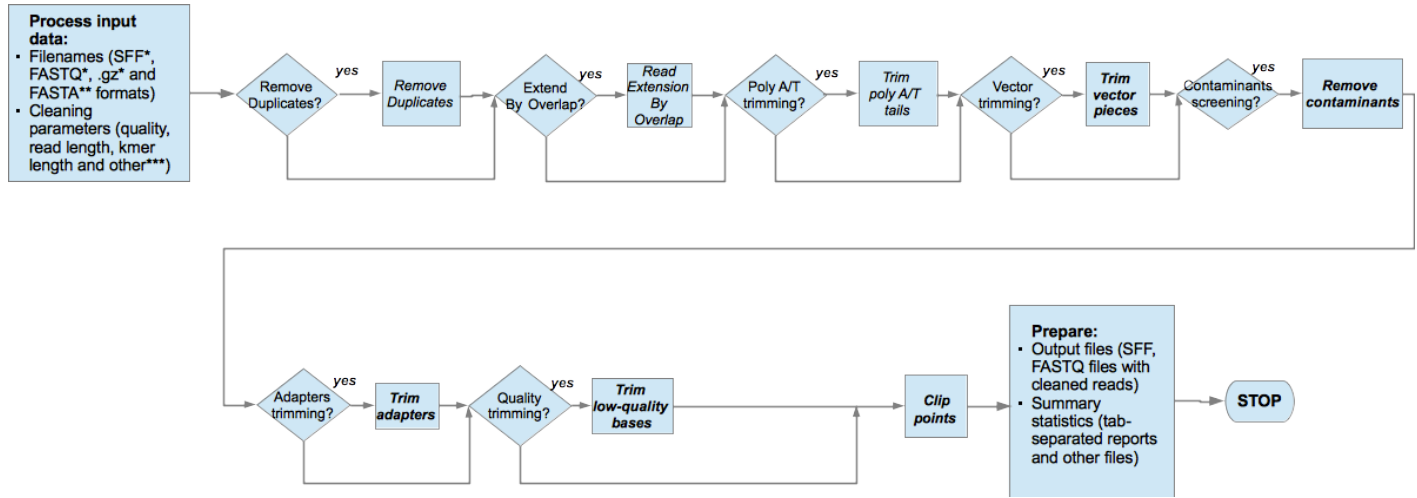


Figure 1: The workflow diagram for SeqyClean.

Table 5: Supported RLMIDs 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	ACGCGTGGTCG
RL4	ACGACACGTAT	ATACGTGGCGT	RL22	CGACGAGTACT	AGTACTGGTCG
RL5	ACGAGTAGACT	AGTCTACGCGT	RL23	CGATACTACGT	ACGTAGTGTCG
RL6	ACGCGTCTAGT	ACTAGAGGCGT	RL24	CGTACGTTCGAT	ATCGACGGACG
RL7	ACGTACACACT	AGTGTGTGCGT	RL25	CTACTCGTAGT	ACTACGGGTAG
RL8	ACGTACTGTGT	ACACAGTGCGT	RL26	GTACAGTACGT	ACGTACGGTAC
RL9	ACGTAGATCGT	ACGATCTGCGT	RL27	GTCGTACGTAT	ATACGTAGGAC
RL10	ACTACGTCTCT	AGAGACGGAGT	RL28	GTGTACGACGT	ACGTCGTGCAC
RL11	ACTATACGAGT	ACTCGTAGAGT	RL29	ACACAGTGAGT	ACTCACGGTGT
RL12	ACTCGCGTCGT	ACGACGGGAGT	RL30	ACACTCATACT	AGTATGGGTGT
RL13	AGACTCGACGT	ACGTCGGGTCT	RL31	ACAGACAGCGT	ACGCTGTGTGT
RL14	AGTACGAGAGT	ACTCTCGGACT	RL32	ACAGACTATAT	ATATAGTGTGT
RL15	AGTACTACTAT	ATAGTAGGACT	RL33	ACAGAGACTCT	AGAGTCTGTGT
RL16	AGTAGACGTCT	AGACGTCGACT	RL34	ACAGCTCGTGT	ACACGAGGTGT
RL17	AGTCGTACACT	AGTGTAGGACT	RL35	ACAGTGTCGAT	ATCGACAGTGT
RL18	AGTGTAGTAGT	ACTACTAGACT	RL36	ACGAGCGCGCT	AGCGCGCGCGT

3 Acknowledgements

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