

Here we provide a full list of our scripts together with the command lines used to benchmark BRAT-nova against Bismark, BS-Seeker2 and BRAT-BW. Please note that the order of arguments in the command lines is important.

Script name	Command Line
Brief description of purpose	Options
<p><i>generate_reads.cpp</i></p> <p>To compile: g++ -o generate_reads generate_reads.cpp</p> <p>Purpose: generates bisulfite-treated single-end reads with sequencing errors, SNPs, indels, adapter contamination</p> <p>Indels and adapter contamination are introduced only if options -i and -a are used</p> <p>Length of indels is random from 1 to 10</p> <p>97% of all Cs are converted to Ts</p>	<p><i>./generate_reads -g references_names.txt -r reads.fastq -p positions.txt -l 100 -t 1000000 -e 5 -s 2 -i 7 -a 10 -m 30</i></p> <p><i>references_names.txt</i> file contains the paths to the files with references (one reference per file): one path per line (option -g)</p> <p><i>reads.fastq</i> contains the resulting reads in FASTQ format (option -r, if not provided, output will be in file <i>reads.fastq</i>)</p> <p><i>positions.txt</i> contains the positions from which the reads have been generated in format: read_id, read, reference_name, strand, position starting with 1, mismatches, indel length, indel's position within the read (option -p, if not provided, output will be in file <i>positions.txt</i>)</p> <p>100 is read length (default 100, option -l)</p> <p>1000000 is the total number of reads generated (default 100, option -t)</p> <p>5 is the percentage of sequencing errors of total number of read bases (default 2, option -e)</p> <p>2 is the percentage of SNPs of total number of reads bases (default 1, option -s)</p> <p>7 is the percentage of reads having an indel at a random position within a read; indel length is random in the range 1...10 (default value for the percentage of reads having an indel is 1, option -i)</p> <p>10 is the percentage of reads with an adapter sequence of random length up to 15bp (next</p>

	<p>argument) replaced the 3'-end of the reads (default is 10%, option -a)</p> <p>30 maximum length of adapter sequence (default is 15bp, option -m)</p>
<p><i>Generate_methylation_call.cpp</i></p> <p>To compile: g++ -o Generate_methylation_call Generate_methylation_call.cpp</p> <p>Purpose: generates bisulfite-treated reads with sequencing errors, SNPs, indels, adapter contamination and records methylation level of each cytosine derived from generated reads</p> <p>Indels and adapter contamination are introduced only if options -i and -a are used</p> <p>97% of all Cs are converted to Ts</p>	<p><i>./Generate_methylation_call -g reference.fa -r reads.fastq -p positions.txt -l 100 -t 1000000 -e 5 -s 2 -i 7 -a 10 -m 30 -o methylation_level.txt</i></p> <p><i>reference.fa</i> is a file with a reference in FASTA format (containing only a single reference, i.e. character ">" occurs only once in the file) (option -g)</p> <p><i>reads.fastq</i> contains the resulting reads in FASTQ format (option -r, if not provided, output will be in file <i>reads.fastq</i>)</p> <p><i>positions.txt</i> contains the positions from which the reads have been generated in format: read_id, read, reference_name, strand, position starting with 1, mismatches, indel length, indel's position within the read (option -p, if not provided, output will be in file <i>positions.txt</i>)</p> <p>100 is read length (default 100, option -l)</p> <p>1000000 is the total number of reads generated (default 100, option -t)</p> <p>5 is the percentage of sequencing errors of total number of read bases (default 2, option -e)</p> <p>2 is the percentage of SNPs of total number of reads bases (default 1, option -s)</p> <p>7 is the percentage of reads having an indel at a random position within a read; indel length equals to 3 (default value for the percentage of reads having an indel is 1, option -i)</p> <p>10 is the percentage of reads with an adapter sequence of random length up to 15bp (next argument) replaced the 3'-end of the reads</p>

	<p>(default is 10%, option -a)</p> <p>30 maximum length of adapter sequence (default is 15bp, option -m)</p> <p><i>methylation_level.txt</i> is a file with methylation level in format: reference name, strand, position starting with 1, count of Cs, count of Cs plus Ts</p>
<p>Correctness_methCall_rand_bratbw.cpp</p> <p>To compile: g++ -o Correctness_methCall_rand_bratbw Correctness_methCall_rand_bratbw.cpp</p> <p>Purpose: it calculates methylation call accuracy and methylation level accuracy, given an output of BRAT-BW' and BRAT-nova' tool <i>acgt_count</i> and given the file <i>methylation_level.txt</i> (produced by running Generate_methylation_call).</p>	<p><i>./Correctness_methCall_rand_bratbw -g reference.fa -m methylation_level.txt -t meth_level_BRAT_nova.txt -r 10 -f 0.2 -c 0.5</i></p> <p><i>reference.fa</i> is the same file used to run Generate_methylation_call (option -g)</p> <p><i>methylation_level.txt</i> is the file produced by running Generate_methylation_call (option -m)</p> <p><i>meth_level_BRAT_nova.txt</i> is the output of <i>acgt_count</i> of BRAT-nova/BRAT-BW (if using BRAT-BW, concatenate two output files) (option -t)</p> <p>10 is the minimum read coverage of a cytosine to consider it for methylation call/level accuracy estimation (default is 10, option -r)</p> <p>0.2 is the threshold on FDR (program will use an error corresponding to 20% of false discovery rate to estimate methylation level accuracy) (default is 0.2, option -f)</p> <p>0.5 is the threshold on methylation level to estimate methylation call accuracy (default is 0.5, option -c)</p>
<p>Correctness_meth_call_BSseeker.cpp</p> <p>To compile: g++ -o Correctness_meth_call_BSseeker Correctness_meth_call_BSseeker.cpp</p> <p>Purpose: it calculates methylation call accuracy and methylation level accuracy, given an output of BS-</p>	<p><i>./Correctness_meth_call_BSseeker -g reference.fa -m methylation_level.txt -t meth_level_BSseeker -r 10 -f 0.2 -c 0.5</i></p> <p>All options are described above (<i>Correctness_methCall_rand_bratbw</i>)</p> <p><i>meth_level_BSseeker</i> is the file produced by</p>

<p>Seeker2's tool <i>bs_seeker2-call_methylation.py</i> (file provided with --CGmap option) and given the file <i>methylation_level.txt</i> (produced by running <i>Generate_methylation_call</i>).</p>	<p><i>bs_seeker2-call_methylation.py</i> (file provided with --CGmap option) in the format: reference name, dummy, pos (starts with 1), dummy, dummy, methylation level, count of Cs, count of Cs and Ts</p> <p>where "dummy" values were ignored.</p>
<p>Correctness_meth_call_bismark.cpp</p> <p>To compile: g++ -o Correctness_meth_call_bismark Correctness_meth_call_bismark.cpp</p> <p>Purpose: it calculates methylation call accuracy and methylation level accuracy, given an output of Bismark's tool <i>bismark_methylation_extractor</i>, the file, concatenation of the output files: cat CpG_O* Non_CpG_O* > meth_level_Bismark.txt</p>	<p><i>./Correctness_meth_call_bismark -g reference.fa -m methylation_level.txt -t meth_level_Bismark.txt -r 10 -f 0.2 -c 0.5</i></p> <p><i>meth_level_Bismark.txt</i> is in format: read id, strand, reference name, position (starts with 1), dummy</p> <p>The rest of parameter options as in <i>Correctness_methCall_rand_bratbw</i></p>
<p>correctness_bratbw.cpp</p> <p>To compile: g++ -o correctness_bratbw correctness_bratbw.cpp</p> <p>Purpose: calculates mapping accuracy of BRAT-BW, given the file <i>positions.txt</i> (results of running <i>generate_reads</i>) and given the output of BRAT-BW</p>	<p><i>./correctness_bratbw -r output_bratbw.txt -p positions.txt -t 1000000 -e 50</i></p> <p>1000000 is the total number of reads produced by <i>generate_reads</i> (attempted to map) (option – t) 50 is the threshold on position, a read is considered to be correctly mapped if it is mapped to the same chromosome, same strand to the position within 50bp of original position (option – e)</p> <p><i>positions.txt</i> is the file with original positions (produced by running <i>generate_reads</i>) (option – p)</p> <p><i>output_bratbw.txt</i> is the output file with mapped reads by BRAT-BW (reads initially generated using <i>generate_reads</i>) (option – r)</p>
<p>correctness_SAM.cpp</p> <p>To compile: g++ -o correctness_SAM correctness_SAM.cpp</p> <p>Purpose: calculates mapping accuracy of Bismark, BS-Seeker2 and BRAT-nova on single-end reads, given the file <i>positions.txt</i> (results of running <i>generate_reads</i>) and given the output of a tool with</p>	<p><i>./correctness_SAM -r mapped_reads.sam -p positions.txt -t 1000000 -e 50</i></p> <p>1000000 is the total number of reads produced by <i>generate_reads</i> (attempted to map) (option – t) 50 is the threshold on position, a read is considered to be correctly mapped if it is mapped to the same chromosome, same strand to the position within 50bp of original position (option – e)</p>

<p>mapped reads in SAM format</p> <p>It also works for paired-end reads for the tools BRAT-nova and Bismark (for Bismark, first run “change_pairs_ids”, described below, and only then run correctness_SAM)</p>	<p>positions.txt is the file with original positions (produced by running generate_reads) (option –p)</p> <p>mapped_reads.sam is the output file with mapped reads (reads initially generated using generate_reads) (option –r)</p>
<p>generate_pairs.cpp</p> <p>To compile: g++ -o generate_pairs generate_pairs.cpp</p> <p>Purpose: generates bisulfite-treated paired-end reads with sequencing errors, SNPs, indels, adapter contamination</p> <p>Indels and adapter contamination are introduced only if options -i and -a are used</p> <p>Length of indels is random from 1 to 10</p> <p>97% of all Cs are converted to Ts</p>	<p>./generate_pairs -g references_names.txt -r paired_end_reads -p positions.txt -l 100 -t 1000000 -e 5 -s 2 -i 7 -a 10 -m 30</p> <p><i>references_names.txt</i> file contains the paths to the files with references (one reference per file): one path per line (option -g)</p> <p>paired-end reads will be in two files: (1) <i>paired_end_reads.mate1.fastq</i> and (2) <i>paired_end_reads.mate2.fastq</i> in FASTQ format (thus, prefix of the file names is provided with option -r, if not provided, prefix is “paired_end_reads”)</p> <p><i>positions.txt</i> contains the positions from which the reads have been generated in format: read_id, read, reference_name, strand, position starting with 1, mismatches, indel length, indel’s position within the read (option -p, if not provided, output will be in file <i>positions.txt</i>)</p> <p>100 is read length (default 100, option -l)</p> <p>1000000 is the total number of reads generated (default 100, option -t)</p> <p>5 is the percentage of sequencing errors of total number of read bases (default 2, option -e)</p> <p>2 is the percentage of SNPs of total number of reads bases (default 1, option -s)</p> <p>7 is the percentage of reads having an indel at a random position within a read; indel length is random in the range 1...10 (default value for the percentage of reads having an indel is 1, option -</p>

	<p>i)</p> <p>10 is the percentage of reads with an adapter sequence of random length up to 15bp (next argument) replaced the 3'-end of the reads (default is 10%, option -a)</p> <p>30 maximum length of adapter sequence (default is 15bp, option -m)</p>
<p>correctness_sam_paired_BSseeker.cpp</p> <p>To compile:</p> <pre>g++ -o correctness_sam_paired_BSseeker correctness_sam_paired_BSseeker.spp</pre> <p>Purpose: calculates mapping accuracy for paired-end reads on BS-Seeker2</p>	<p><i>./correctness_sam_paired_BSseeker -r mapped_reads.sam -p positions.txt -t 1000000 -e 50</i></p> <p>All options are the same as described in correctness_SAM</p> <p>1000000 here is the total number of attempted reads (count single reads, i.e. if there were 500,000 pairs, then we need to pass 1000000 = 2*500000 with option -t)</p>
<p>change_pairs_ids.cpp</p> <p>To compile:</p> <pre>g++ -o change_pairs_ids change_pairs_ids.cpp</pre> <p>Purpose: before calculating of mapping accuracy for paired-end reads on Bismark, run this script first, then use correctness_SAM</p>	<p><i>./change_pairs_ids output_paired_Bismark.sam</i></p> <p>output_paired_Bismark.sam is the output of Bismark mapping paired-end reads</p> <p>The output with adjusted reads ids will be in <i>output_paired_Bismark.sam.out</i> (use this file to calculate mapping accuracy of Bismark on paired-end reads)</p>
<p>correctness_bratbw_paired.cpp</p> <p>To compile:</p> <pre>g++ -o correctness_bratbw_paired correctness_bratbw_paired.cpp</pre> <p>Purpose: calculates mapping accuracy of BRAT-bw on paired-end reads</p>	<p><i>./correctness_bratbw_paired -r output_bratbw.txt -p positions.txt -t 1000000 -e 50</i></p> <p>1000000 is the total number of reads produced by generate_pairs (attempted to map, total number of pairs multiplied by 2) (option -t)</p> <p>50 is the threshold on position, a read is considered to be correctly mapped if it is mapped to the same chromosome, same strand to the position within 50bp of original position (option -e)</p> <p>positions.txt is the file with original positions</p>

	<p>(produced by running <code>generate_pairs</code>) (option <code>-p</code>)</p> <p><code>output_bratbw.txt</code> is the output file with mapped reads by BRAT-BW (reads initially generated using <code>generate_reads</code>) (option <code>-r</code>)</p>
To convert FASTQ reads to BRAT-BW's input reads format, use tool <i>trim</i> in BRAT-BW suit.	