

---

# **viral-ngs Documentation**

***Release v2.0.13***

**Broad Institute Viral Genomics**

**2019-11-07**



<b>1</b>	<b>Contents</b>	<b>1</b>
1.1	Description of the methods . . . . .	1
1.2	Command line tools . . . . .	1



---

## Contents

---

### 1.1 Description of the methods

This is a base module that provides basic utility functions, some short read aligners, and samtools and Picard.

### 1.2 Command line tools

#### 1.2.1 reports.py - produce various metrics and reports

Functions to create reports from genomics pipeline data.

```
usage: reports.py subcommand
```

##### Sub-commands:

##### assembly\_stats

Fetch assembly-level statistics for a given sample

```
usage: reports.py assembly_stats [-h]
                                [--cov_thresholds COV_THRESHOLDS [COV_
                                ↪THRESHOLDS ...]]
                                [--assembly_dir ASSEMBLY_DIR]
                                [--assembly_tmp ASSEMBLY_TMP]
                                [--align_dir ALIGN_DIR]
                                [--reads_dir READS_DIR]
                                [--raw_reads_dir RAW_READS_DIR]
                                [--loglevel
                                ↪{DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                [--version] [--tmp_dir TMP_DIR]
                                [--tmp_dirKeep]
                                samples [samples ...] outFile
```

##### Positional arguments:

<b>samples</b>	Sample names.
<b>outFile</b>	Output report file.

##### Options:

**--cov\_thresholds=(1, 5, 20, 100)** Genome coverage thresholds to report on. (default: %(default)s)

**--assembly\_dir=data/02\_assembly** Directory with assembly outputs. (default: %(default)s)

**--assembly\_tmp=tmp/02\_assembly** Directory with assembly temp files. (default: %(default)s)

**--align\_dir=data/02\_align\_to\_self** Directory with reads aligned to own assembly. (default: %(default)s)

**--reads\_dir=data/01\_per\_sample** Directory with unaligned filtered read BAMs. (default: %(default)s)

**--raw\_reads\_dir=data/00\_raw** Directory with unaligned raw read BAMs. (default: %(default)s)

**--loglevel=INFO** Verboseness of output. [default: %(default)s]  
Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION

**--version, -V** show program's version number and exit

**--tmp\_dir=/tmp** Base directory for temp files. [default: %(default)s]

**--tmp\_dirKeep=False** Keep the tmp\_dir if an exception occurs while running. Default is to delete all temp files at the end, even if there's a failure.

#### coverage\_only

```
usage: reports.py coverage_only [-h]
                                [--cov_thresholds COV_THRESHOLDS [COV_
                                ↪THRESHOLDS ...]]
                                [--loglevel
                                ↪{DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                [--version] [--tmp_dir TMP_DIR]
                                [--tmp_dirKeep]
                                mapped_bams [mapped_bams ...] out_report
```

#### Positional arguments:

**mapped\_bams** Aligned-to-self mapped bam files.

**out\_report** Output report file.

#### Options:

**--cov\_thresholds=(1, 5, 20, 100)** Genome coverage thresholds to report on. (default: %(default)s)

**--loglevel=INFO** Verboseness of output. [default: %(default)s]  
Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION

**--version, -V** show program's version number and exit

**--tmp\_dir=/tmp** Base directory for temp files. [default: %(default)s]

**--tmp\_dirKeep=False** Keep the tmp\_dir if an exception occurs while running. Default is to delete all temp files at the end, even if there's a failure.

**alignment\_summary**

Write or print pairwise alignment summary information for sequences in two FASTA files, including SNPs, ambiguous bases, and indels.

```
usage: reports.py alignment_summary [-h] [--outfileName OUTFILENAME]
                                   [--printCounts]
                                   [--loglevel
    ↪ {DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                   [--version] [--tmp_dir TMP_DIR]
                                   [--tmp_dirKeep]
                                   inFastaFileOne inFastaFileTwo
```

**Positional arguments:**

<b>inFastaFileOne</b>	First fasta file for an alignment
<b>inFastaFileTwo</b>	First fasta file for an alignment

**Options:**

<b>--outfileName</b>	Output file for counts in TSV format
<b>--printCounts=False</b>	Undocumented
<b>--loglevel=INFO</b>	Verboseness of output. [default: %(default)s] Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION
<b>--version, -V</b>	show program's version number and exit
<b>--tmp_dir=/tmp</b>	Base directory for temp files. [default: %(default)s]
<b>--tmp_dirKeep=False</b>	Keep the tmp_dir if an exception occurs while running. Default is to delete all temp files at the end, even if there's a failure.

**consolidate\_fastqc**

Consolidate multiple FASTQC reports into one.

```
usage: reports.py consolidate_fastqc [-h]
                                     [--loglevel
    ↪ {DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                     [--version] [--tmp_dir TMP_DIR]
                                     [--tmp_dirKeep]
                                     inDirs [inDirs ...] outFile
```

**Positional arguments:**

<b>inDirs</b>	Input FASTQC directories.
<b>outFile</b>	Output report file.

**Options:**

<b>--loglevel=INFO</b>	Verboseness of output. [default: %(default)s] Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION
<b>--version, -V</b>	show program's version number and exit
<b>--tmp_dir=/tmp</b>	Base directory for temp files. [default: %(default)s]

**--tmp\_dirKeep=False** Keep the tmp\_dir if an exception occurs while running. Default is to delete all temp files at the end, even if there's a failure.

### consolidate\_spike\_count

Consolidate multiple spike count reports into one.

```
usage: reports.py consolidate_spike_count [-h]
                                         [--loglevel
                                         ↪ {DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                         [--version] [--tmp_dir TMP_DIR]
                                         [--tmp_dirKeep]
                                         inDir outFile
```

#### Positional arguments:

<b>in_dir</b>	Input spike count directory.
<b>out_file</b>	Output report file.

#### Options:

<b>--loglevel=INFO</b>	Verboseness of output. [default: %(default)s] Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION
<b>--version, -V</b>	show program's version number and exit
<b>--tmp_dir=/tmp</b>	Base directory for temp files. [default: %(default)s]
<b>--tmp_dirKeep=False</b>	Keep the tmp_dir if an exception occurs while running. Default is to delete all temp files at the end, even if there's a failure.

### aggregate\_spike\_count

aggregate multiple spike count reports into one.

```
usage: reports.py aggregate_spike_count [-h]
                                         [--loglevel
                                         ↪ {DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                         [--version] [--tmp_dir TMP_DIR]
                                         [--tmp_dirKeep]
                                         inDir outFile
```

#### Positional arguments:

<b>in_dir</b>	Input spike count directory.
<b>out_file</b>	Output report file.

#### Options:

<b>--loglevel=INFO</b>	Verboseness of output. [default: %(default)s] Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION
<b>--version, -V</b>	show program's version number and exit
<b>--tmp_dir=/tmp</b>	Base directory for temp files. [default: %(default)s]
<b>--tmp_dirKeep=False</b>	Keep the tmp_dir if an exception occurs while running. Default is to delete all temp files at the end, even if there's a failure.



**plot\_coverage**

Generate a coverage plot from an aligned bam file

```
usage: reports.py plot_coverage [-h] [--plotFormat] [--plotDataStyle]
                                [--plotStyle] [--plotWidth PLOT_WIDTH]
                                [--plotHeight PLOT_HEIGHT]
                                [--plotDPI PLOT_DPI] [--plotTitle PLOT_TITLE]
                                [--plotXLimits PLOT_X_LIMITS PLOT_X_LIMITS]
                                [--plotYLimits PLOT_Y_LIMITS PLOT_Y_LIMITS]
                                [-q BASE_Q_THRESHOLD] [-Q MAPPING_Q_THRESHOLD]
                                [-m MAX_COVERAGE_DEPTH]
                                [-l READ_LENGTH_THRESHOLD] [--binLargePlots]
                                [--binningSummaryStatistic {max,min}]
                                [--outSummary OUT_SUMMARY]
                                [--plotOnlyNonDuplications]
                                [--loglevel
                                ↪ {DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                [--version] [--tmp_dir TMP_DIR]
                                [--tmp_dirKeep]
                                in_bam out_plot_file
```

**Positional arguments:**

**in\_bam** Input reads, BAM format.

**out\_plot\_file** The generated chart file

**Options:**

**--plotFormat** File format of the coverage plot. By default it is inferred from the file extension of out\_plot\_file, but it can be set explicitly via --plotFormat. Valid formats include: ps, eps, pdf, pgf, png, raw, rgba, svg, svgz, jpg, jpeg, tif, tiff

Possible choices: ps, eps, pdf, pgf, png, raw, rgba, svg, svgz, jpg, jpeg, tif, tiff

**--plotDataStyle=filled** The plot data display style. Valid options: filled, line, dots (default: %(default)s)

Possible choices: filled, line, dots

**--plotStyle=ggplot** The plot visual style. Valid options: seaborn-white, seaborn-pastel, seaborn-deep, seaborn-darkgrid, dark\_background, seaborn-paper, grayscale, seaborn-muted, tableau-colorblind10, fivethirtyeight, seaborn-poster, seaborn, fast, seaborn-dark, seaborn-whitegrid, bmh, seaborn-bright, seaborn-dark-palette, \_classic\_test, ggplot, seaborn-notebook, seaborn-colorblind, Solarize\_Light2, classic, seaborn-talk, seaborn-ticks (default: %(default)s)

Possible choices: seaborn-white, seaborn-pastel, seaborn-deep, seaborn-darkgrid, dark\_background, seaborn-paper, grayscale, seaborn-muted, tableau-colorblind10, fivethirtyeight, seaborn-poster, seaborn, fast, seaborn-dark, seaborn-whitegrid, bmh, seaborn-bright, seaborn-dark-palette, \_classic\_test, ggplot, seaborn-notebook, seaborn-colorblind, Solarize\_Light2, classic, seaborn-talk, seaborn-ticks

**--plotWidth=880** Width of the plot in pixels (default: %(default)s)

**--plotHeight=680** Width of the plot in pixels (default: %(default)s)

**--plotDPI=100.0** dots per inch for rendered output, more useful for vector modes (default: %(default)s)

**--plotTitle=Coverage Plot** The title displayed on the coverage plot (default: ‘%(default)s’)

**--plotXLimits** Limits on the x-axis of the coverage plot; args are ‘<min> <max>’

**--plotYLimits** Limits on the y-axis of the coverage plot; args are ‘<min> <max>’

**-q** The minimum base quality threshold

**-Q** The minimum mapping quality threshold

**-m** The max coverage depth (default: %(default)s)

**-l** Read length threshold

**--binLargePlots=False** Plot summary read depth in one-pixel-width bins for large plots.

**--binningSummaryStatistic=max** Statistic used to summarize each bin (max or min).  
Possible choices: max, min

**--outSummary** Coverage summary TSV file. Default is to write to temp.

**--plotOnlyNonDuplicates=False** Plot only non-duplicates (samtools -F 1024), coverage counted by bedtools rather than samtools.

**--loglevel=INFO** Verboseness of output. [default: %(default)s]  
Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION

**--version, -V** show program’s version number and exit

**--tmp\_dir=/tmp** Base directory for temp files. [default: %(default)s]

**--tmp\_dirKeep=False** Keep the tmp\_dir if an exception occurs while running. Default is to delete all temp files at the end, even if there’s a failure.

### align\_and\_plot\_coverage

Take reads, align to reference with BWA-MEM, and generate a coverage plot

```
usage: reports.py align_and_plot_coverage [-h] [--plotFormat]
                                         [--plotDataStyle] [--plotStyle]
                                         [--plotWidth PLOT_WIDTH]
                                         [--plotHeight PLOT_HEIGHT]
                                         [--plotDPI PLOT_DPI]
                                         [--plotTitle PLOT_TITLE]
                                         [--plotXLimits PLOT_X_LIMITS PLOT_X_
↪LIMITS]
                                         [--plotYLimits PLOT_Y_LIMITS PLOT_Y_
↪LIMITS]
                                         [-q BASE_Q_THRESHOLD]
                                         [-Q MAPPING_Q_THRESHOLD]
                                         [-m MAX_COVERAGE_DEPTH]
                                         [-l READ_LENGTH_THRESHOLD]
                                         [--binLargePlots]
```

```

↪]
                                [--binningSummaryStatistic {max,min}

                                [--outSummary OUT_SUMMARY]
                                [--outBam OUT_BAM] [--sensitive]
                                [--excludeDuplicates]
                                [--JVMmemory JVMMEMORY]
                                [--picardOptions [PICARDOPTIONS_
↪[PICARDOPTIONS ...]]]
                                [--minScoreToFilter MIN_SCORE_TO_
↪FILTER]
                                [--aligner {novoalign,bwa}]
                                [--aligner_options ALIGNER_OPTIONS]
                                [--NOVOALIGN_LICENSE_PATH NOVOALIGN_
↪LICENSE_PATH]
                                [--loglevel
↪{DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                [--version] [--tmp_dir TMP_DIR]
                                [--tmp_dirKeep]
                                in_bam out_plot_file ref_fasta

```

**Positional arguments:**

<b>in_bam</b>	Input reads, BAM format.
<b>out_plot_file</b>	The generated chart file
<b>ref_fasta</b>	Reference genome, FASTA format.

**Options:**

<b>--plotFormat</b>	File format of the coverage plot. By default it is inferred from the file extension of out_plot_file, but it can be set explicitly via --plotFormat. Valid formats include: ps, eps, pdf, pgf, png, raw, rgba, svg, svgz, jpg, jpeg, tif, tiff  Possible choices: ps, eps, pdf, pgf, png, raw, rgba, svg, svgz, jpg, jpeg, tif, tiff
<b>--plotDataStyle=filled</b>	The plot data display style. Valid options: filled, line, dots (default: %(default)s)  Possible choices: filled, line, dots
<b>--plotStyle=ggplot</b>	The plot visual style. Valid options: seaborn-white, seaborn-pastel, seaborn-deep, seaborn-darkgrid, dark_background, seaborn-paper, grayscale, seaborn-muted, tableau-colorblind10, fivethirtyeight, seaborn-poster, seaborn, fast, seaborn-dark, seaborn-whitegrid, bmh, seaborn-bright, seaborn-dark-palette, _classic_test, ggplot, seaborn-notebook, seaborn-colorblind, Solarize_Light2, classic, seaborn-talk, seaborn-ticks (default: %(default)s)  Possible choices: seaborn-white, seaborn-pastel, seaborn-deep, seaborn-darkgrid, dark_background, seaborn-paper, grayscale, seaborn-muted, tableau-colorblind10, fivethirtyeight, seaborn-poster, seaborn, fast, seaborn-dark, seaborn-whitegrid, bmh, seaborn-bright, seaborn-dark-palette, _classic_test, ggplot, seaborn-notebook, seaborn-colorblind, Solarize_Light2, classic, seaborn-talk, seaborn-ticks

<b>--plotWidth=880</b>	Width of the plot in pixels (default: <code>%(default)s</code> )
<b>--plotHeight=680</b>	Width of the plot in pixels (default: <code>%(default)s</code> )
<b>--plotDPI=100.0</b>	dots per inch for rendered output, more useful for vector modes (default: <code>%(default)s</code> )
<b>--plotTitle=Coverage Plot</b>	The title displayed on the coverage plot (default: <code>'%(default)s'</code> )
<b>--plotXLimits</b>	Limits on the x-axis of the coverage plot; args are <code>'&lt;min&gt; &lt;max&gt;'</code>
<b>--plotYLimits</b>	Limits on the y-axis of the coverage plot; args are <code>'&lt;min&gt; &lt;max&gt;'</code>
<b>-q</b>	The minimum base quality threshold
<b>-Q</b>	The minimum mapping quality threshold
<b>-m</b>	The max coverage depth (default: <code>%(default)s</code> )
<b>-l</b>	Read length threshold
<b>--binLargePlots=False</b>	Plot summary read depth in one-pixel-width bins for large plots.
<b>--binningSummaryStatistic=max</b>	Statistic used to summarize each bin (max or min). Possible choices: max, min
<b>--outSummary</b>	Coverage summary TSV file. Default is to write to temp.
<b>--outBam</b>	Output aligned, indexed BAM file. Default is to write to temp.
<b>--sensitive=False</b>	Equivalent to giving bwa: <code>'-k 12 -A 1 -B 1 -O 1 -E 1'</code> . Only relevant if the bwa aligner is selected (the default).
<b>--excludeDuplicates=False</b>	MarkDuplicates with Picard and only plot non-duplicates
<b>--JVMmemory=2g</b>	JVM virtual memory size (default: <code>%(default)s</code> )
<b>--picardOptions=[]</b>	Optional arguments to Picard's MarkDuplicates, <code>OPTION-NAME=value ...</code>
<b>--minScoreToFilter</b>	Filter bwa alignments using this value as the minimum allowed alignment score. Specifically, sum the alignment scores across all alignments for each query (including reads in a pair, supplementary and secondary alignments) and then only include, in the output, queries whose summed alignment score is at least this value. This is only applied when the aligner is 'bwa'. The filtering on a summed alignment score is sensible for reads in a pair and supplementary alignments, but may not be reasonable if bwa outputs secondary alignments (i.e., if '-a' is in the aligner options). (default: not set - i.e., do not filter bwa's output)
<b>--aligner=bwa</b>	aligner (default: <code>%(default)s</code> ) Possible choices: novoalign, bwa
<b>--aligner_options</b>	aligner options (default for novoalign: <code>"-r Random -l 40 -g 40 -x 20 -t 100 -k"</code> , bwa: bwa defaults)

**--NOVOALIGN\_LICENSE\_PATH** A path to the novoalign.lic file. This overrides the NOVOALIGN\_LICENSE\_PATH environment variable. (default: %(default)s)

**--loglevel=INFO** Verboseness of output. [default: %(default)s]  
Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION

**--version, -V** show program's version number and exit

**--tmp\_dir=/tmp** Base directory for temp files. [default: %(default)s]

**--tmp\_dirKeep=False** Keep the tmp\_dir if an exception occurs while running. Default is to delete all temp files at the end, even if there's a failure.

**fastqc**

```
usage: reports.py fastqc [-h] inBam outHtml
```

**Positional arguments:**

<b>inBam</b>	Input reads, BAM format.
<b>outHtml</b>	Output report, HTML format.

**1.2.2 illumina.py - for raw Illumina outputs**

Utilities for demultiplexing Illumina data.

```
usage: illumina.py subcommand
```

**Sub-commands:****illumina\_demux**

Read Illumina runs & produce BAM files, demultiplexing to one bam per sample, or for simplex runs, a single bam will be produced bearing the flowcell ID. Wraps together Picard's ExtractBarcodes (for multiplexed samples) and IlluminaBasecallsToSam while handling the various required input formats. Also can read Illumina BCL directories, tar.gz BCL directories.

```
usage: illumina.py illumina_demux [-h] [--outMetrics OUTMETRICS]
                                   [--commonBarcodes COMMONBARCODES]
                                   [--sampleSheet SAMPLESHEET]
                                   [--runInfo RUNINFO] [--flowcell FLOWCELL]
                                   [--read_structure READ_STRUCTURE]
                                   [--max_mismatches MAX_MISMATCHES]
                                   [--minimum_base_quality MINIMUM_BASE_
↪QUALITY]
                                   [--min_mismatch_delta MIN_MISMATCH_DELTA]
                                   [--max_no_calls MAX_NO_CALLS]
                                   [--minimum_quality MINIMUM_QUALITY]
                                   [--compress_outputs COMPRESS_OUTPUTS]
                                   [--sequencing_center SEQUENCING_CENTER]
                                   [--adapters_to_check [ADAPTERS_TO_CHECK_
↪[ADAPTERS_TO_CHECK ...]]]
                                   [--platform PLATFORM]
                                   [--max_reads_in_ram_per_tile MAX_READS_IN_
↪RAM_PER_TILE]
                                   [--max_records_in_ram MAX_RECORDS_IN_RAM]
```

```

[--apply_eamss_filter APPLY_EAMSS_FILTER]
[--force_gc FORCE_GC]
[--first_tile FIRST_TILE]
[--tile_limit TILE_LIMIT]
[--include_non_pf_reads INCLUDE_NON_PF_]
↪ READS]

[--run_start_date RUN_START_DATE]
[--read_group_id READ_GROUP_ID]
[--compression_level COMPRESSION_LEVEL]
[--JVMmemory JVMMEMORY] [--threads THREADS]
[--loglevel
↪ {DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
[--version] [--tmp_dir TMP_DIR]
[--tmp_dirKeep]
inDir lane outDir

```

**Positional arguments:**

<b>inDir</b>	Illumina BCL directory (or tar.gz of BCL directory). This is the top-level run directory.
<b>lane</b>	Lane number.
<b>outDir</b>	Output directory for BAM files.

**Options:**

<b>--outMetrics</b>	Output ExtractIlluminaBarcodes metrics file. Default is to dump to a temp file.
<b>--commonBarcodes</b>	Write a TSV report of all barcode counts, in descending order. Only applicable for read structures containing “B”
<b>--sampleSheet</b>	Override SampleSheet. Input tab or CSV file w/header and four named columns: barcode_name, library_name, barcode_sequence_1, barcode_sequence_2. Default is to look for a SampleSheet.csv in the inDir.
<b>--runInfo</b>	Override RunInfo. Input xml file. Default is to look for a RunInfo.xml file in the inDir.
<b>--flowcell</b>	Override flowcell ID (default: read from RunInfo.xml).
<b>--read_structure</b>	Override read structure (default: read from RunInfo.xml).
<b>--max_mismatches=0</b>	Picard ExtractIlluminaBarcodes MAX_MISMATCHES (default: %(default)s)
<b>--minimum_base_quality=20</b>	Picard ExtractIlluminaBarcodes MINIMUM_BASE_QUALITY (default: %(default)s)
<b>--min_mismatch_delta</b>	Picard ExtractIlluminaBarcodes MIN_MISMATCH_DELTA (default: %(default)s)
<b>--max_no_calls</b>	Picard ExtractIlluminaBarcodes MAX_NO_CALLS (default: %(default)s)
<b>--minimum_quality</b>	Picard ExtractIlluminaBarcodes MINIMUM_QUALITY (default: %(default)s)
<b>--compress_outputs</b>	Picard ExtractIlluminaBarcodes COMPRESS_OUTPUTS (default: %(default)s)

**--sequencing\_center** Picard IlluminaBasecallsToSam SEQUENCING\_CENTER (default: %(default)s)

**--adapters\_to\_check=('PAIRED\_END', 'NEXTERA\_V1', 'NEXTERA\_V2')** Picard IlluminaBasecallsToSam ADAPTERS\_TO\_CHECK (default: %(default)s)

**--platform** Picard IlluminaBasecallsToSam PLATFORM (default: %(default)s)

**--max\_reads\_in\_ram\_per\_tile=1000000** Picard IlluminaBasecallsToSam MAX\_READS\_IN\_RAM\_PER\_TILE (default: %(default)s)

**--max\_records\_in\_ram=2000000** Picard IlluminaBasecallsToSam MAX\_RECORDS\_IN\_RAM (default: %(default)s)

**--apply\_eamss\_filter** Picard IlluminaBasecallsToSam APPLY\_EAMSS\_FILTER (default: %(default)s)

**--force\_gc** Picard IlluminaBasecallsToSam FORCE\_GC (default: %(default)s)

**--first\_tile** Picard IlluminaBasecallsToSam FIRST\_TILE (default: %(default)s)

**--tile\_limit** Picard IlluminaBasecallsToSam TILE\_LIMIT (default: %(default)s)

**--include\_non\_pf\_reads=False** Picard IlluminaBasecallsToSam INCLUDE\_NON\_PF\_READS (default: %(default)s)

**--run\_start\_date** Picard IlluminaBasecallsToSam RUN\_START\_DATE (default: %(default)s)

**--read\_group\_id** Picard IlluminaBasecallsToSam READ\_GROUP\_ID (default: %(default)s)

**--compression\_level=7** Picard IlluminaBasecallsToSam COMPRESSION\_LEVEL (default: %(default)s)

**--JVMmemory=7g** JVM virtual memory size (default: %(default)s)

**--threads=0** Number of threads (default: 0)

**--loglevel=INFO** Verboseness of output. [default: %(default)s]  
Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION

**--version, -V** show program's version number and exit

**--tmp\_dir=/tmp** Base directory for temp files. [default: %(default)s]

**--tmp\_dirKeep=False** Keep the tmp\_dir if an exception occurs while running. Default is to delete all temp files at the end, even if there's a failure.

### lane\_metrics

Write out lane metrics to a tsv file.

```
usage: illumina.py lane_metrics [-h] [--read_structure READ_STRUCTURE]
                                [--JVMmemory JVMMEMORY]
                                [--loglevel
                                ↪ {DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                [--version] [--tmp_dir TMP_DIR]
```

```

[--tmp_dirKeep]
inDir outPrefix

```

**Positional arguments:**

<b>inDir</b>	Illumina BCL directory (or tar.gz of BCL directory). This is the top-level run directory.
<b>outPrefix</b>	Prefix path to the *.illumina_lane_metrics and *.illumina_phasing_metrics files.

**Options:**

<b>--read_structure</b>	Override read structure (default: read from RunInfo.xml).
<b>--JVMmemory=8g</b>	JVM virtual memory size (default: %(default)s)
<b>--loglevel=INFO</b>	Verboseness of output. [default: %(default)s]  Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION
<b>--version, -V</b>	show program's version number and exit
<b>--tmp_dir=/tmp</b>	Base directory for temp files. [default: %(default)s]
<b>--tmp_dirKeep=False</b>	Keep the tmp_dir if an exception occurs while running. Default is to delete all temp files at the end, even if there's a failure.

**common\_barcode**s

Extract Illumina barcodes for a run and write a TSV report of the barcode counts in descending order

```

usage: illumina.py common_barcode [-h] [--truncateToLength TRUNCATETOLENGTH]
                                  [--omitHeader] [--includeNoise]
                                  [--outMetrics OUTMETRICS]
                                  [--sampleSheet SAMPLESHEET]
                                  [--flowcell FLOWCELL]
                                  [--read_structure READ_STRUCTURE]
                                  [--max_mismatches MAX_MISMATCHES]
                                  [--minimum_base_quality MINIMUM_BASE_
↳QUALITY]
                                  [--min_mismatch_delta MIN_MISMATCH_DELTA]
                                  [--max_no_calls MAX_NO_CALLS]
                                  [--minimum_quality MINIMUM_QUALITY]
                                  [--compress_outputs COMPRESS_OUTPUTS]
                                  [--JVMmemory JVMMEMORY]
                                  [--loglevel
↳{DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                  [--version] [--tmp_dir TMP_DIR]
                                  [--tmp_dirKeep]
                                  inDir lane outSummary

```

**Positional arguments:**

<b>inDir</b>	Illumina BCL directory (or tar.gz of BCL directory). This is the top-level run directory.
<b>lane</b>	Lane number.
<b>outSummary</b>	Path to the summary file (.tsv format). It includes several columns: (barcode1, likely_index_name1, barcode2, likely_index_name2, count), where likely index names are either



the exact match index name for the barcode sequence, or those Hamming distance of 1 away.

#### Options:

- truncateToLength** If specified, only this number of barcodes will be returned. Useful if you only want the top N barcodes.
- omitHeader=False** If specified, a header will not be added to the outSummary tsv file.
- includeNoise=False** If specified, barcodes with periods (".") will be included.
- outMetrics** Output ExtractIlluminaBarcodes metrics file. Default is to dump to a temp file.
- sampleSheet** Override SampleSheet. Input tab or CSV file w/header and four named columns: barcode\_name, library\_name, barcode\_sequence\_1, barcode\_sequence\_2. Default is to look for a SampleSheet.csv in the inDir.
- flowcell** Override flowcell ID (default: read from RunInfo.xml).
- read\_structure** Override read structure (default: read from RunInfo.xml).
- max\_mismatches=0** Picard ExtractIlluminaBarcodes MAX\_MISMATCHES (default: %(default)s)
- minimum\_base\_quality=20** Picard ExtractIlluminaBarcodes MINIMUM\_BASE\_QUALITY (default: %(default)s)
- min\_mismatch\_delta** Picard ExtractIlluminaBarcodes MIN\_MISMATCH\_DELTA (default: %(default)s)
- max\_no\_calls** Picard ExtractIlluminaBarcodes MAX\_NO\_CALLS (default: %(default)s)
- minimum\_quality** Picard ExtractIlluminaBarcodes MINIMUM\_QUALITY (default: %(default)s)
- compress\_outputs** Picard ExtractIlluminaBarcodes COMPRESS\_OUTPUTS (default: %(default)s)
- JVMmemory=8g** JVM virtual memory size (default: %(default)s)
- loglevel=INFO** Verboseness of output. [default: %(default)s]  
Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION
- version, -V** show program's version number and exit
- tmp\_dir=/tmp** Base directory for temp files. [default: %(default)s]
- tmp\_dirKeep=False** Keep the tmp\_dir if an exception occurs while running. Default is to delete all temp files at the end, even if there's a failure.

#### guess\_barcodes

Guess the barcode value for a sample name, based on the following: - a list is made of novel barcode pairs seen in the data, but not in the picard metrics - for the sample in question, get the most abundant novel barcode pair where one of the barcodes seen in the data matches one of the barcodes in the picard metrics (partial match) - if there are no partial matches, get the most abundant novel barcode pair

Limitations: - If multiple samples share a barcode with multiple novel barcodes, disentangling them is difficult or impossible

The names of samples to guess are selected: - explicitly by name, passed via argument, OR - explicitly by read count threshold, OR - automatically (if names or count threshold are omitted) based on basic outlier detection of deviation from an assumed-balanced pool with some number of negative controls

```
usage: illumina.py guess_barcodes [-h]
                                [--readcount_threshold READCOUNT_THRESHOLD_]
↪ | --sample_names [SAMPLE_NAMES [SAMPLE_NAMES ...]]]
                                [--outlier_threshold OUTLIER_THRESHOLD]
                                [--expected_assigned_fraction EXPECTED_
↪ ASSIGNED_FRACTION]
                                [--number_of_negative_controls NUMBER_OF_
↪ NEGATIVE_CONTROLS]
                                [--rows_limit ROWS_LIMIT]
                                [--loglevel
↪ {DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                [--version] [--tmp_dir TMP_DIR]
                                [--tmp_dirKeep]
                                in_barcodes in_picard_metrics
                                out_summary_tsv
```

#### Positional arguments:

<b>in_barcodes</b>	The barcode counts file produced by common_barcodes.
<b>in_picard_metrics</b>	The demultiplexing read metrics produced by Picard.
<b>out_summary_tsv</b>	Path to the summary file (.tsv format). It includes several columns: (sample_name, expected_barcode_1, expected_barcode_2, expected_barcode_1_name, expected_barcode_2_name, expected_barcodes_read_count, guessed_barcode_1, guessed_barcode_2, guessed_barcode_1_name, guessed_barcode_2_name, guessed_barcodes_read_count, match_type), where the expected values are those used by Picard during demultiplexing and the guessed values are based on the barcodes seen among the data.

#### Options:

<b>--readcount_threshold</b>	If specified, guess barcodes for samples with fewer than this many reads.
<b>--sample_names</b>	If specified, only guess barcodes for these sample names.
<b>--outlier_threshold=0.675</b>	threshold of how far from unbalanced a sample must be to be considered an outlier.
<b>--expected_assigned_fraction=0.7</b>	The fraction of reads expected to be assigned. An exception is raised if fewer than this fraction are assigned.
<b>--number_of_negative_controls=1</b>	The number of negative controls in the pool, for calculating expected number of reads in the rest of the pool.
<b>--rows_limit=1000</b>	The number of rows to use from the in_barcodes.
<b>--loglevel=INFO</b>	Verboseness of output. [default: %(default)s]  Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION

**--version, -V** show program's version number and exit

**--tmp\_dir=/tmp** Base directory for temp files. [default: %(default)s]

**--tmp\_dirKeep=False** Keep the tmp\_dir if an exception occurs while running. Default is to delete all temp files at the end, even if there's a failure.

### miseq\_fastq\_to\_bam

Convert fastq read files to a single bam file. Fastq file names must conform to patterns emitted by Miseq machines. Sample metadata must be provided in a SampleSheet.csv that corresponds to the fastq filename. Specifically, the `_S##_` index in the fastq file name will be used to find the corresponding row in the SampleSheet

```
usage: illumina.py miseq_fastq_to_bam [-h] [--inFastq2 INFASTQ2]
                                     [--runInfo RUNINFO]
                                     [--sequencing_center SEQUENCING_CENTER]
                                     [--JVMmemory JVMMEMORY]
                                     [--loglevel
                                     ↪ {DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                     [--version] [--tmp_dir TMP_DIR]
                                     [--tmp_dirKeep]
                                     outBam sampleSheet inFastq1
```

#### Positional arguments:

**outBam** Output BAM file.

**sampleSheet** Input SampleSheet.csv file.

**inFastq1** Input fastq file; 1st end of paired-end reads if paired.

#### Options:

**--inFastq2** Input fastq file; 2nd end of paired-end reads.

**--runInfo** Input RunInfo.xml file.

**--sequencing\_center** Name of your sequencing center (default is the sequencing machine ID from the RunInfo.xml)

**--JVMmemory=2g** JVM virtual memory size (default: %(default)s)

**--loglevel=INFO** Verboseness of output. [default: %(default)s]  
Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION

**--version, -V** show program's version number and exit

**--tmp\_dir=/tmp** Base directory for temp files. [default: %(default)s]

**--tmp\_dirKeep=False** Keep the tmp\_dir if an exception occurs while running. Default is to delete all temp files at the end, even if there's a failure.

### extract\_fc\_metadata

Extract RunInfo.xml and SampleSheet.csv from the provided Illumina directory

```
usage: illumina.py extract_fc_metadata [-h]
                                     [--loglevel
                                     ↪ {DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                     [--version] [--tmp_dir TMP_DIR]
                                     [--tmp_dirKeep]
                                     flowcell outRunInfo outSampleSheet
```

**Positional arguments:**

<b>flowcell</b>	Illumina directory (possibly tarball)
<b>outRunInfo</b>	Output RunInfo.xml file.
<b>outSampleSheet</b>	Output SampleSheet.csv file.

**Options:**

<b>--loglevel=INFO</b>	Verboseness of output. [default: %(default)s] Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION
<b>--version, -V</b>	show program's version number and exit
<b>--tmp_dir=/tmp</b>	Base directory for temp files. [default: %(default)s]
<b>--tmp_dirKeep=False</b>	Keep the tmp_dir if an exception occurs while running. Default is to delete all temp files at the end, even if there's a failure.

### 1.2.3 broad\_utils.py - for data generated at the Broad Institute

Utilities for getting sequences out of the Broad walk-up sequencing pipeline. These utilities are probably not of much use outside the Broad.

```
usage: broad_utils.py subcommand
```

**Sub-commands:****get\_bustard\_dir**

Find the basecalls directory from a Picard directory

```
usage: broad_utils.py get_bustard_dir [-h]
                                     [--loglevel
                                     ↪ {DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                     inDir
```

**Positional arguments:**

<b>inDir</b>	Picard directory
--------------	------------------

**Options:**

<b>--loglevel=ERROR</b>	Verboseness of output. [default: %(default)s] Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION
-------------------------	---

**get\_run\_date**

Find the sequencing run date from a Picard directory

```
usage: broad_utils.py get_run_date [-h]
                                   [--loglevel
                                   ↪ {DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                   inDir
```

**Positional arguments:**

<b>inDir</b>	Picard directory
--------------	------------------

**Options:**

**--loglevel=ERROR** Verboseness of output. [default: %(default)s]  
 Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION

**get\_all\_names**

Get all samples

```
usage: broad_utils.py get_all_names [-h]
                                   [--loglevel
                                   ↪ {DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                   {samples, libraries, runs} runfile
```

**Positional arguments:**

**type** Type of name  
 Possible choices: samples, libraries, runs

**runfile** File with seq run information

**Options:**

**--loglevel=ERROR** Verboseness of output. [default: %(default)s]  
 Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION

## 1.2.4 file\_utils.py - utilities to perform various file manipulations

Utilities for dealing with files.

```
usage: file_utils.py subcommand
```

**Sub-commands:****merge\_tarballs**

Merges separate tarballs into one tarball data can be piped in and/or out

```
usage: file_utils.py merge_tarballs [-h]
                                   [--extractToDiskPath EXTRACT_TO_DISK_PATH]
                                   [--pipeInHint PIPE_HINT_IN]
                                   [--pipeOutHint PIPE_HINT_OUT]
                                   [--threads THREADS]
                                   [--loglevel
                                   ↪ {DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION}]
                                   [--version] [--tmp_dir TMP_DIR]
                                   [--tmp_dirKeep]
                                   out_tarball in_tarballs [in_tarballs ...]
```

**Positional arguments:**

**out\_tarball** output tarball (\*.tar.gz|\*.tar.lz4|\*.tar.bz2|\*.tar.zstl-); compression is inferred by the file extension. Note: if “-” is used, output will be written to stdout and --pipeOutHint must be provided to indicate compression type when compression type is not gzip (gzip is used by default).

**in\_tarballs** input tarballs (\*.tar.gz|\*.tar.lz4|\*.tar.bz2|\*.tar.zst)

**Options:**

**--extractToDiskPath** If specified, the tar contents will also be extracted to a local directory.

**--pipeInHint=gz** If specified, the compression type used is used for piped input.

**--pipeOutHint=gz** If specified, the compression type used is used for piped output.

**--threads** Number of threads (default: all available cores)

**--loglevel=INFO** Verboseness of output. [default: %(default)s]  
Possible choices: DEBUG, INFO, WARNING, ERROR, CRITICAL, EXCEPTION

**--version, -V** show program's version number and exit

**--tmp\_dir=/tmp** Base directory for temp files. [default: %(default)s]

**--tmp\_dirKeep=False** Keep the tmp\_dir if an exception occurs while running. Default is to delete all temp files at the end, even if there's a failure.