

---

## Supplementary information

---

# Capture-C: a modular and flexible approach for high-resolution chromosome conformation capture

---

In the format provided by the authors and unedited

---

## Supplementary information

---

# Capture-C: a modular and flexible approach for high-resolution chromosome conformation capture

---

In the format provided by the authors and unedited

---

**ccanalyser**

***Release 0.0.1***

**asmith, dsims**

**Apr 15, 2021**



# CONTENTS

<b>1</b>	<b>Installing</b>	<b>3</b>
1.1	Pre-Installation recommendations . . . . .	3
1.2	Installation . . . . .	3
<b>2</b>	<b>Pipeline</b>	<b>5</b>
2.1	Step 1 - Create a working directory . . . . .	5
2.2	Step 2 - Edit a copy of config.yml . . . . .	5
2.3	Step 3 - Copy or link fastq files into the working directory . . . . .	6
2.4	Step 4 - Running the pipeline . . . . .	6
2.5	Step 5 - Running the pipeline to a specified stage . . . . .	7
<b>3</b>	<b>Pipeline FAQ</b>	<b>9</b>
3.1	Interruptions to the pipeline . . . . .	9
3.2	DRMAA error . . . . .	9
3.3	KeyError: ‘missing parameter accessed’ . . . . .	9
3.4	ValueError: not enough values to unpack (expected 2, got 1) . . . . .	10
<b>4</b>	<b>CLI Documentation</b>	<b>11</b>
4.1	ccanalyser . . . . .	11
<b>5</b>	<b>Module Documentation</b>	<b>25</b>
5.1	CCanalyser Modules . . . . .	25
<b>6</b>	<b>Glossary</b>	<b>63</b>
	<b>Python Module Index</b>	<b>65</b>
	<b>Index</b>	<b>67</b>



CCanalyser is a package explicitly designed for processing Capture-C, Tri-C and Tiled-C data. It consists of a configurable data processing pipeline and a supporting command line interface to enable fine-grained control over the analysis. Unlike other pipelines that are designed to process Hi-C or Capture-HiC data, the filtering steps in ccanalyser are specifically optimised Capture-C, Tri-C and Tiled-C. This package is effectively a python based re-implementation of [CCseqBasicS](#), [Tri-C](#), [Tiled-C](#) and [CaptureCompare](#) designed to provide an end-to-end solution for processing data from these three assays.



## INSTALLING

### 1.1 Pre-Installation recommendations

1. Install conda if it has not been already using the [conda install instructions](#).
2. If ccanalyser is **not** being installed through conda, first generate a new conda environment using the yaml file in the [GitHub repository](#):

```
conda env create -f ccanalyser_conda_env.yml  
conda activate cc
```

3. If you intend to use a cluster e.g. SunGrid engine/SLURM add the path to the DRMAA interface to your .bashrc:

```
# Access to the DRMAA library: https://en.wikipedia.org/wiki/DRMAA  
echo "export DRMAA_LIBRARY_PATH=/<full-path>/libdrmaa.so" >> ~/.bashrc  
  
# You can get this value from your configured environment:  
env | grep DRMAA_LIBRARY_PATH  
  
# or just look for the library:  
find / -name "*libdrmaa.so"
```

### 1.2 Installation

The package can be installed in several ways:

---

**Note:** Currently only github installation is supported

---

1. Install from conda:

```
conda install ccanalyser # do not use: not available yet
```

2. Install from pypi:

```
pip install ccanalyser # do not use: not available yet
```

3. Install from GitHub:

```
git clone https://github.com/sims-lab/capture-c.git
cd capture-c
pip install .
```

**PIPELINE**

The main feature of ccanalyser is the end-to-end data processing pipeline. The pipeline has been written using the [cgat-core workflow management system](#) and the following diagram illustrates the steps performed by the pipeline:

This section provides further details on how to run the pipeline. In essence the pipeline requires a working directory with correctly named FASTQ files and a [config.yml](#) file that provides the pipeline configuration.

## 2.1 Step 1 - Create a working directory

To run the pipeline you will need to create a working directory for the pipeline run:

```
mkdir RS411_EPZ5676/
cd RS411_EPZ5676/
```

The pipeline will be executed here and all files will be generated in this directory.

## 2.2 Step 2 - Edit a copy of config.yml

The configuration file [config.yml](#) enables parameterisation of the pipeline run with user specific settings. Furthermore, it also provides paths to essential files for the pipeline run (e.g., bowtie2 indices). The paths supplied do not have to be in the same directory as the pipeline but it is recommended to copy the capture viewpoints used to the working directory.

**Warning:** The yaml file must be named **config.yml** for the pipeline to recognise it and run correctly.

This yaml file can be edited using standard text editors e.g.:

```
# To use gedit
gedit config.yml

# To use nano
nano config.yml
```

## 2.3 Step 3 - Copy or link fastq files into the working directory

The pipeline requires that fastq files are paired and in any of these formats:

---

**Note:** Gzipped files are handled appropriately without the need for extraction if .gz is present at the end of the file name.

---

**Note:** Multi-lane FASTQ files should be concatenated prior to running the pipeline; otherwise multiple separate analyses will be performed.

---

Here is an example of file pairing for two samples:

- samplename1\_R1.fastq.gz
- samplename1\_R2.fastq.gz
- samplename2\_1.fastq
- samplename2\_2.fastq

All FASTQ files present in the directory will be processed by the pipeline in parallel and original FASTQ files will not be modified. If new FASTQ files are added to a pre-run pipeline, only the new files will be processed.

Copy:

```
cp PATH_TO_FASTQ/example_R1.fastq.gz.
```

Symlink example:

Be sure to use the absolute path for symlinks

```
ln -s /ABSOLUTE_PATH_TO_FASTQ/example_R1.fastq.gz
```

## 2.4 Step 4 - Running the pipeline

After copying/linking FASTQ files into the working directory and configuring the copy of `config.yml` in the working directory for the current experiment, the pipeline can be run with:

```
ccanalyser pipeline
```

There are several options to visualise which tasks will be performed by the pipeline before running.

The tasks to be performed can be examined with:

```
# Shows the tasks to be performed
ccanalyser pipeline show

# Plots a directed graph using graphviz
ccanalyser pipeline plot
```

If you are happy with the tasks to be performed, the full pipeline run can be launched with:

```
# If using all default settings and using a cluster
ccanalyser pipeline make

# If not using a cluster, run in local mode.
ccanalyser pipeline make --local -p 4

# Avoiding network disconnections
nohup ccanalyser pipeline make &
```

See [cgat-core Read the Docs](#) for additional information.

## 2.5 Step 5 - Running the pipeline to a specified stage

There are currently multiple stopping points built into the pipeline at key stages. These are:

- fastq\_preprocessing - Stops after *in silico* digestion of FASTQ files.
- pre\_annotation - Stops before aligned slices are ready to be annotated.
- post\_annotation - Stops after aligned slices have been annotated.
- post\_ccanalyser\_analysis - Stops after reporters have been identified and duplicate filtered.
- full - Run the pipeline until all required tasks are complete.

To run the pipeline until one of these stopping points, use:

```
# Run until TASK_NAME step
ccanalyser pipeline make TASK_NAME
```



## PIPELINE FAQ

### 3.1 Interruptions to the pipeline

If for any reason the pipeline interrupted, simply delete any malformed file or directory and re-run the pipeline. The pipeline will then continue with the analysis without needing to start from the beginning.

Leftover ctmpXXX.sh files can be safely deleted if the pipeline has stopped running:

```
rm -f *.sh*
```

To restart the pipeline from the beginning for any reason, simply delete all processed files and re-run the pipeline:

```
rm -rf ccanalyser_preprocessing/ ccanalyser_analysis/ ccanalyser_compare/ statistics/_  
↳pipeline.log  
ccanalyser pipeline make
```

### 3.2 DRMAA error

A common error when running the pipeline is:

```
GLOBAL_SESSION = drmaa.Session()  
NameError: name 'drmaa' is not defined
```

This error occurs because you are not connected to the cluster. See [Pre-Installation recommendations](#) for how to add DRMAA\_LIBRARY\_PATH to your bashrc. Alternatively, if you do not want to use a cluster you can run the pipeline in *local mode*.

### 3.3 KeyError: ‘missing parameter accessed’

This occurs if one of the keys in config.yml has been altered or deleted. Either find the source of the error or get a new version of the yaml file and fill it out.

### 3.4 ValueError: not enough values to unpack (expected 2, got 1)

The fastq files are not paired correctly. Please ensure that fastq files are in the format:

```
NAME-OF-SAMPLE_(R)1.fastq.gz  
NAME-OF-SAMPLE_(R)2.fastq.gz
```

## CLI DOCUMENTATION

### 4.1 ccanalyser

An end to end solution for processing: Capture-C, Tri-C and Tiled-C data.

```
ccanalyser [OPTIONS] COMMAND [ARGS]...
```

#### 4.1.1 alignments

Alignment annotation, identification and deduplication.

```
ccanalyser alignments [OPTIONS] COMMAND [ARGS]...
```

##### annotate

Annotates a bed file with other bed files using bedtools intersect.

Whilst bedtools intersect allows for interval names and counts to be used for annotating intervals, this command provides the ability to annotate intervals with both interval names and counts at the same time. As the pipeline allows for empty bed files, this command has built in support to deal with blank/malformed bed files and will return default N/A values.

Prior to interval annotation, the bed file to be intersected is validated and duplicate entries/multimapping reads are removed to ensure consistent annotations and prevent issues with reporter identification.

```
ccanalyser alignments annotate [OPTIONS] SLICES
```

##### Options

**-a, --actions <actions>**

Determines if the overlaps are counted or if the name should just be reported

**Options** get | count

**-b, --bed\_files <bed\_files>**

Bed file(s) to intersect with slices

**-n, --names <names>**

Names to use as column names for the output tsv file.

```
-f, --overlap_fractions <overlap_fractions>
    The minimum overlap required for an intersection between two intervals to be reported.

-o, --output <output>
    Path for the annotated slices to be output.

--duplicates <duplicates>
    Method to use for reconciling duplicate slices (i.e. multimapping). Currently only ‘remove’ is supported.

        Options remove

-p, --n_cores <n_cores>
    Intersections are performed by chromosome, this determines the number of cores.

--invalid_bed_action <invalid_bed_action>
    Method to deal with invalid bed files e.g. blank or incorrectly formatted. Setting this to ‘ignore’ will report
    default N/A values (either ‘.’ or 0) for invalid files

        Options ignore | error
```

## Arguments

### SLICES

Required argument

### deduplicate

Identifies and removes duplicated aligned fragments.

PCR duplicates are very commonly present in Capture-C/Tri-C/Tiled-C data and must be removed for accurate analysis. Unlike fastq deduplicate, this command removes fragments with identical genomic coordinates.

Non-combined (pe) and combined (flashed) reads are treated slightly differently due to the increased confidence that the ligation junction has been captured for the flashed reads.

```
ccanalyser alignments deduplicate [OPTIONS] COMMAND [ARGS] ...
```

### identify

```
ccanalyser alignments deduplicate identify [OPTIONS] FRAGMENTS_FN
```

## Options

```
-o, --output <output>
    Path for outputting fragments with duplicated coordinates in json format.

--buffer <buffer>
    Number of fragments to process at one time in order to preserve memory.

--read_type <read_type>
    Indicates if the fragments have been combined (flashed) or not (pe).

        Options flashed | pe
```

## Arguments

### **FRAGMENTS\_FN**

Required argument

### **remove**

Removes duplicated aligned fragments.

Parses a tsv file containing aligned read slices and outputs only slices from unique fragments. Duplicated parental read id determined by the “identify” subcommand are located within the slices tsv file and removed.

Outputs statistics for the number of unique slices and the number of duplicate slices identified.

```
ccanalyser alignments deduplicate remove [OPTIONS] SLICES_FN
```

## Options

### **-d, --duplicated\_ids <duplicated\_ids>**

Path to duplicated fragment ids determined by the ‘identify’ subcommand.

### **-o, --output <output>**

Path for outputting deduplicated slices in tsv format.

### **--buffer <buffer>**

Number of fragments to process at one time, in order to preserve memory.

### **--stats\_prefix <stats\_prefix>**

Output prefix for deduplication statistics

### **--sample\_name <sample\_name>**

Name of sample being analysed e.g. DOX\_treated\_1. Required for correct statistics.

### **--read\_type <read\_type>**

Indicates if the fragments have been combined (flashed) or not (pe). Required for correct statistics.

**Options** flashed | pe

## Arguments

### **SLICES\_FN**

Required argument

### **filter**

Removes unwanted aligned slices and identifies reporters.

Parses a BAM file and merges this with a supplied annotation to identify unwanted slices. Filtering can be tuned for Capture-C, Tri-C and Tiled-C data to ensure optimal filtering.

```
ccanalyser alignments filter [OPTIONS] [capture|tri|tiled]
```

## Options

**-b, --bam <bam>**  
**Required** Bam file to process

**-a, --annotations <annotations>**  
**Required** Annotations for the bam file that must contain the required columns, see description.

**--custom\_filtering <custom\_filtering>**  
Custom filtering to be used. This must be supplied as a path to a yaml file.

**-o, --output\_prefix <output\_prefix>**  
Output prefix for deduplicated fastq file(s)

**--stats\_prefix <stats\_prefix>**  
Output prefix for stats file(s)

**--sample\_name <sample\_name>**  
Name of sample e.g. DOX\_treated\_1

**--read\_type <read\_type>**  
Type of read  
    **Options** flashed | pe

**--gzip, --no-gzip**  
Determines if files are gzipped or not

## Arguments

### METHOD

Required argument

## 4.1.2 fastq

Fastq splitting, deduplication and digestion.

```
ccanalyser fastq [OPTIONS] COMMAND [ARGS] ...
```

### deduplicate

Identifies PCR duplicate fragments from Fastq files.

PCR duplicates are very commonly present in Capture-C/Tri-C/Tiled-C data and must be removed for accurate analysis. These commands attempt to identify and remove duplicate reads/fragments from fastq file(s) to speed up downstream analysis.

```
ccanalyser fastq deduplicate [OPTIONS] COMMAND [ARGS] ...
```

## identify

Identifies fragments with duplicated sequences.

Merges the hashed dictionaries (in json format) generated by the “parse” subcommand and identifies read with exactly the same sequence (share an identical hash). Duplicated read identifiers (hashed) are output in json format. The “remove” subcommand uses this dictionary to remove duplicates from fastq files.

```
ccanalyser fastq deduplicate identify [OPTIONS] [INPUT_FILES]...
```

### Options

**-o, --output <output>**  
Required Output file

### Arguments

#### INPUT\_FILES

Optional argument(s)

## parse

Parses fastq file(s) into easy to deduplicate format.

This command parses one or more fastq files and generates a dictionary containing hashed read identifiers together with hashed concatenated sequences. The hash dictionary is output in json format and the identify subcommand can be used to determine which read identifiers have duplicate sequences.

```
ccanalyser fastq deduplicate parse [OPTIONS] INPUT_FILES...
```

### Options

**-o, --output <output>**  
File to store hashed sequence identifiers  
**--read\_buffer <read\_buffer>**  
Number of reads to process before writing to file

### Arguments

#### INPUT\_FILES

Required argument(s)

## **remove**

Removes fragments with duplicated sequences from fastq files.

Parses input fastq files and removes any duplicates from the fastq file(s) that are present in the json file supplied. This json dictionary should be produced by the “identify” subcommand.

Statistics for the number of duplicated and unique reads are also provided.

```
ccanalyser fastq deduplicate remove [OPTIONS] [INPUT_FILES]...
```

## **Options**

**-o, --output\_prefix <output\_prefix>**  
Output prefix for deduplicated fastq file(s)

**-d, --duplicated\_ids <duplicated\_ids>**  
Path to duplicate ids, identified by the identify subcommand

**--read\_buffer <read\_buffer>**  
Number of reads to process before writing to file

**--gzip, --no-gzip**  
Determines if files are gziped or not

**--compression\_level <compression\_level>**  
Level of compression for output files

**--sample\_name <sample\_name>**  
Name of sample e.g. DOX\_treated\_1

**--stats\_prefix <stats\_prefix>**  
Output prefix for stats file

## **Arguments**

**INPUT\_FILES**  
Optional argument(s)

## **digest**

Performs in silico digestion of one or a pair of fastq files.

```
ccanalyser fastq digest [OPTIONS] INPUT_FASTQ...
```

## Options

**-r, --restriction\_enzyme** <restriction\_enzyme>  
**Required** Restriction enzyme name or sequence to use for in silico digestion.

**-m, --mode** <mode>  
**Required** Digestion mode. Combined (Flashed) or non-combined (PE) read pairs.  
**Options** flashed | pe

**-o, --output\_file** <output\_file>

**-p, --n\_cores** <n\_cores>

**--minimum\_slice\_length** <minimum\_slice\_length>

**--keep\_cutsite** <keep\_cutsite>

**--compression\_level** <compression\_level>  
Level of compression for output files (1=low, 9=high)

**--read\_buffer** <read\_buffer>  
Number of reads to process before writing to file to conserve memory.

**--stats\_prefix** <stats\_prefix>  
Output prefix for stats file

**--sample\_name** <sample\_name>  
Name of sample e.g. DOX\_treated\_1. Required for correct statistics.

## Arguments

**INPUT\_FASTQ**  
Required argument(s)

## split

Splits fastq file(s) into equal chunks of n reads.

```
ccanalyser fastq split [OPTIONS] INPUT_FILES...
```

## Options

**-m, --method** <method>  
Method to use for splitting  
**Options** python | unix

**-o, --output\_prefix** <output\_prefix>  
Output prefix for deduplicated fastq file(s)

**--compression\_level** <compression\_level>  
Level of compression for output files

**-n, --n\_reads** <n\_reads>  
Number of reads per fastq file

**--gzip, --no-gzip**  
Determines if files are gzipped or not

## Arguments

### **INPUT\_FILES**

Required argument(s)

## 4.1.3 genome

Genome wide methods digestion.

```
ccanalyser genome [OPTIONS] COMMAND [ARGS]...
```

### digest

Performs in silico digestion of a genome in fasta format.

Digests the supplied genome fasta file and generates a bed file containing the locations of all restriction fragments produced by the supplied restriction enzyme.

A log file recording the number of restriction fragments for the supplied genome is also generated.

```
ccanalyser genome digest [OPTIONS] INPUT_FASTA
```

## Options

**-r, --recognition\_site <recognition\_site>**

Required Recognition enzyme or sequence

**-l, --logfile <logfile>**

Path for digestion log file

**-o, --output\_file <output\_file>**

Output file path

**--remove\_cutsite <remove\_cutsite>**

Exclude the recognition sequence from the output

**--sort**

Sorts the output bed file by chromosome and start coord.

## Arguments

### **INPUT\_FASTA**

Required argument

#### 4.1.4 pipeline

Runs the data processing pipeline

```
ccanalyser pipeline [OPTIONS] [make|show|clone|touch] [PIPELINE_OPTIONS]...
```

##### Options

**-h, --help**

##### Arguments

###### MODE

Required argument

###### PIPELINE\_OPTIONS

Optional argument(s)

#### 4.1.5 reporters

Reporter counting, storing, comparison, pileups and heatmaps.

```
ccanalyser reporters [OPTIONS] COMMAND [ARGS]...
```

##### count

Determines the number of captured restriction fragment interactions genome wide.

Parses a reporter slices tsv and counts the number of unique restriction fragment interaction combinations that occur within each fragment.

Options to ignore unwanted counts e.g. excluded regions or capture fragments are provided. In addition the number of reporter fragments can be subsampled if required.

```
ccanalyser reporters count [OPTIONS] REPORTERS
```

##### Options

**-o, --output <output>**

Name of output file

**--remove\_exclusions**

Prevents analysis of fragments marked as proximity exclusions

**--remove\_capture**

Prevents analysis of capture fragment interactions

**--subsample <subsample>**

Subsamples reporters before analysis of interactions

## Arguments

### REPORTERS

Required argument

### differential

Identifies differential interactions between conditions.

Parses a union bedgraph containing reporter counts from at least two conditions with two or more replicates for a single capture probe and outputs differential interaction results. Following filtering to ensure that the number of interactions is above the required threshold (`-threshold_count`), `diffxpy` is used to run a wald test after fitting a negative binomial model to the interaction counts. The options to filter results can be filtered by a minimum mean value (`threshold_mean`) and/or maximum q-value (`threshold_q`) are also provided.

Notes:

Currently both the capture viewpoints and the name of the probe being analysed must be provided in order to correctly extract cis interactions.

If a `N_SAMPLE * METADATA` design matrix has not been supplied, the script assumes that the standard replicate naming structure has been followed i.e. `SAMPLE_CONDITION_REPLICATE_(1|2).fastq.gz`.

```
ccanalyser reporters differential [OPTIONS] UNION_BEDGRAPH
```

## Options

- `-n, --capture_name <capture_name>`**  
**Required** Name of capture probe, must be present in viewpoint file.
- `-c, --capture_viewpoints <capture_viewpoints>`**  
**Required** Path to capture viewpoints bed file
- `-o, --output_prefix <output_prefix>`**  
Output prefix for pairwise statistical comparisons
- `--design_matrix <design_matrix>`**  
Path tsv file containing sample annotations (`N_SAMPLES * N_INFO_COLUMNS`)
- `--grouping_col <grouping_col>`**  
Column to use for grouping replicates
- `--threshold_count <threshold_count>`**  
Minimum count required to be considered for analysis
- `--threshold_q <threshold_q>`**  
Upper threshold of q-value required for output.
- `--threshold_mean <threshold_mean>`**  
Minimum mean count required for output.

## Arguments

### **UNION\_BEDGRAPH**

Required argument

## pileup

Extracts reporters from a capture experiment and generates a bedgraph file.

Identifies reporters for a single probe (if a probe name is supplied) or all capture probes present in a capture experiment HDF5 file.

The bedgraph generated can be normalised by the number of cis interactions for inter experiment comparisons and/or binned into even genomic windows.

```
ccanalyser reporters pileup [OPTIONS] COOLER_FN
```

## Options

### **-n, --capture\_names <capture\_names>**

Capture to extract and convert to bedgraph, if not provided will transform all.

### **-o, --output\_prefix <output\_prefix>**

Output prefix for bedgraphs

### **--normalise**

Normalised bedgraph (Correct for number of cis reads)

### **--binsize <binsize>**

Binsize to use for converting bedgraph to evenly sized genomic bins

### **--gzip**

Compress output using gzip

### **--scale\_factor <scale\_factor>**

Scale factor to use for bedgraph normalisation

### **--sparse, --dense**

Produce bedgraph containing just positive bins (sparse) or all bins (dense)

## Arguments

### **COOLER\_FN**

Required argument

## plot

Plots a heatmap of reporter interactions.

Parses a HDF5 file containing the result of a capture experiment (binned into even genomic windows) and plots a heatmap of interactions over a specified genomic range. If a capture probe name is not supplied the script will plot all probes present in the file.

Heatmaps can also be normalised (**--normalise**) using either:

- n\_interactions: The number of cis interactions.

- **n\_rf\_n\_interactions:** Normalised to the number of restriction fragments making up both genomic bins and by the number of cis interactions.
- **ice:** ICE normalisation followed by number of cis interactions correction.

```
ccanalyser reporters plot [OPTIONS] COOLER_FN
```

## Options

**-c, --coordinates <coordinates>**  
Coordinates in the format chr1:1000-2000 or a path to a .bed file with coordinates

**-r, --resolution <resolution>**  
**Required** Resolution at which to plot. Must be present within the cooler file.

**-n, --capture\_names <capture\_names>**  
Capture names to plot. If not supplied will plot all

**--normalisation <normalisation>**  
Normalisation method for heatmap  
    **Options** n\_interactions | n\_rf\_n\_interactions | ice

**--cmap <cmap>**  
Colour map to use for plotting

**--vmax <vmax>**  
Vmax for plotting

**--vmin <vmin>**  
Vmin for plotting

**-o, --output\_prefix <output\_prefix>**  
Output prefix for plot

**--remove\_capture**  
Removes the capture probe bins from the matrix

## Arguments

**COOLER\_FN**  
Required argument

## store

Store reporter counts.

These commands store and manipulate reporter restriction fragment interaction counts as cooler formatted groups in HDF5 files.

See subcommands for details.

```
ccanalyser reporters store [OPTIONS] COMMAND [ARGS] ...
```

## bins

Convert a cooler group containing restriction fragments to constant genomic windows

Parses a cooler group and aggregates restriction fragment interaction counts into genomic bins of a specified size. If the normalise option is selected, columns containing normalised counts are added to the pixels table of the output

```
ccanalyser reporters store bins [OPTIONS] COOLER_FN
```

## Options

**-b, --binsizes <binsizes>**  
Binsizes to use for windowing

**--normalise**  
Enables normalisation of interaction counts during windowing

**--overlap\_fraction <overlap\_fraction>**  
Minimum overlap between genomic bins and restriction fragments for overlap

**-p, --n\_cores <n\_cores>**  
Number of cores used for binning

**--scale\_factor <scale\_factor>**  
Scaling factor used for normalisation

**--conversion\_tables <conversion\_tables>**  
Pickle file containing pre-computed fragment -> bin conversions.

**-o, --output <output>**  
Name of output file. (Cooler formatted hdf5 file)

## Arguments

**COOLER\_FN**  
Required argument

## fragments

Stores restriction fragment interaction combinations at the restriction fragment level.

Parses reporter restriction fragment interaction counts produced by “ccanalyser reporters count” and generates a cooler formatted group in an HDF5 File. See <https://cooler.readthedocs.io/en/latest/> for further details.

```
ccanalyser reporters store fragments [OPTIONS] COUNTS
```

## Options

**-f, --fragment\_map** <fragment\_map>  
Required Path to digested genome bed file

**-c, --capture\_viewpoints** <capture\_viewpoints>  
Required Path to capture viewpoints file

**-n, --capture\_name** <capture\_name>  
Required Name of capture viewpoint to store

**-g, --genome** <genome>  
Name of genome

**--suffix** <suffix>  
Suffix to append after the capture name for the output file

**-o, --output** <output>  
Name of output file. (Cooler formatted hdf5 file)

## Arguments

**COUNTS**  
Required argument

### merge

Merges ccanalyser cooler files together.

Produces a unified cooler with both restriction fragment and genomic bins whilst reducing the storage space required by hard linking the “bins” tables to prevent duplication.

```
ccanalyser reporters store merge [OPTIONS] COOLERS...
```

## Options

**-o, --output** <output>  
Output file name

## Arguments

**COOLERS**  
Required argument(s)

---

## MODULE DOCUMENTATION

---

### 5.1 CCanalyser Modules

#### 5.1.1 CCanalyser CLI Modules

##### alignments annotate

```
ccanalyser.cli.alignments_annotate.cycle_argument(arg)
```

Allows for the same argument to be stated once but repeated for all files

```
ccanalyser.cli.alignments_annotate.remove_duplicates_from_bed(bed: Union[str,  
pybed-  
tools.bedtool.BedTool,  
pan-  
das.core.frame.DataFrame])  
→ pybed-  
tools.bedtool.BedTool
```

Simple removal of duplicated entries from bed file.

If a “score” field is present a higher scored entry is prioritised.

**Parameters** `bed`(`Union[str, BedTool, pd.DataFrame]`) – Bed object to deduplicate

**Returns** BedTool with deduplicated names

**Return type** BedTool

```
ccanalyser.cli.alignments_annotate.annotate(slices: os.PathLike, actions: Optional[Tuple]  
= None, bed_files: Optional[Tuple] = None,  
names: Optional[Tuple] = None, over-  
lap_fractions: Optional[Tuple] = None, out-  
put: Optional[os.PathLike] = None, dupli-  
cates: str = 'remove', n_cores: int = 8, in-  
valid_bed_action: str = 'error')
```

Annotates a bed file with other bed files using bedtools intersect.

Whilst bedtools intersect allows for interval names and counts to be used for annotating intervals, this command provides the ability to annotate intervals with both interval names and counts at the same time. As the pipeline allows for empty bed files, this command has built in support to deal with blank/malformed bed files and will return default N/A values.

Prior to interval annotation, the bed file to be intersected is validated and duplicate entries/multimapping reads are removed to ensure consistent annotations and prevent issues with reporter identification.

**Parameters**

- **slices** (*os.PathLike*) – Input bed file.
- **actions** (*Tuple, optional*) – Methods to use for annotation. Choose from (getCount). Defaults to None.
- **bed\_files** (*Tuple, optional*) – Bed files to intersect with the bed file to be annotated. Defaults to None.
- **names** (*Tuple, optional*) – Column names for output tsv file. Defaults to None.
- **overlap\_fractions** (*Tuple, optional*) – Minimum overlap fractions required to call an intersection. Defaults to None.
- **output** (*os.PathLike, optional*) – Output file path for annotated .tsv file. Defaults to None.
- **duplicates** (*str, optional*) – Method to deal with multimapping reads/duplicate bed names. Currently, “remove” is the only supported option. Defaults to “remove”.
- **n\_cores** (*int, optional*) – Number of cores to use for intersection. Bed files are split by chromosome for faster intersection. Defaults to 8.
- **invalid\_bed\_action** (*str, optional*) – Action to deal with invalid bed files. Choose from (ignoreerror) .These can be ignored by setting to “ignore”. Defaults to ‘error’.

**Raises** `NotImplementedError` – Only supported option for duplicate bed names is remove.

## alignments deduplicate

```
ccanalyser.cli.alignments_deduplicate.identify(fragments_fn: os.PathLike, output: os.PathLike = 'duplicated_ids.json', buffer: int = 1000000.0, read_type: str = 'flashed')
```

Identifies aligned fragments with duplicate coordinates.

Parses a tsv file containing filtered aligned fragments and generates a dictionary containing the hashed parental read id and hashed genomic coordinates of all slices. Duplicated fragments are implicitly removed if they share the same genomic coordinate hash.

For non-combined reads (pe) a genomic coordinate hash is generated from the start of the first slice and the end of the last slice. This is due to the decreased confidence in the quality of the centre of the fragment. The coordinate hash for combined reads (flashed) is generated directly from the fragment coordinates. Only fragments with the exact coordinates and slice order will be considered to be duplicates.

Identified duplicate fragments are output in json format to be used by the “remove” subcommand.

### Parameters

- **fragments\_fn** (*os.PathLike*) – Input fragments.tsv file to process.
- **output** (*os.PathLike, optional*) – Output path to output duplicated parental read ids. Defaults to “duplicated\_ids.json”.
- **buffer** (*int, optional*) – Number of fragments to process in memory. Defaults to 1e6.
- **read\_type** (*str, optional*) – Process combined(flashed) or non-combined reads (pe). Due to the low confidence in the quality of pe reads, duplicates are identified by removing any fragments with matching start and end coordinates. Defaults to “flashed”.

```
ccanalyser.cli.alignments_deduplicate.remove(slices_fn: os.PathLike, duplicated_ids: os.PathLike, output: os.PathLike = 'dedup.slices.tsv.gz', buffer: int = 5000000.0, sample_name: str = "", read_type: str = "", stats_prefix: os.PathLike = "")
```

Removes duplicated aligned fragments.

Parses a tsv file containing aligned read slices and outputs only slices from unique fragments. Duplicated parental read id determined by the “identify” subcommand are located within the slices tsv file and removed.

Outputs statistics for the number of unique slices and the number of duplicate slices identified.

#### Parameters

- **slices\_fn** (*os.PathLike*) – Input slices.tsv file.
- **duplicated\_ids** (*os.PathLike*) – Duplicated parental read ids in json format.
- **output** (*os.PathLike, optional*) – Output file path for deduplicated slices. Defaults to “dedup.slices.tsv.gz”.
- **buffer** (*int, optional*) – Number of slices to process in memory. Defaults to 1e6.
- **sample\_name** (*str, optional*) – Name of sample being processed e.g. DOX-treated\_1 used for statistics. Defaults to “”.
- **read\_type** (*str, optional*) – Process combined(flashed) or non-combined reads (pe) used for statistics. Defaults to “”.
- **stats\_prefix** (*os.PathLike, optional*) – Output path for deduplication statistics. Defaults to “”.

### alignments filter

```
ccanalyser.cli.alignments_filter.merge_annotations(df: pd.core.frame.DataFrame, annotations: os.PathLike) → pandas.core.frame.DataFrame
```

Combines annotations with the parsed bam file output.

Uses pandas outer join on the indexes to merge annotations e.g. number of capture probe overlaps.

Annotation tsv must have the index as the first column and this index must have intersecting keys with the first dataframe’s index.

#### Parameters

- **df** (*pd.DataFrame*) – Dataframe to merge with annotations
- **annotations** (*os.PathLike*) – Filename of .tsv to read and merge with df

**Returns** Merged dataframe

**Return type** *pd.DataFrame*

```
ccanalyser.cli.alignments_filter.filter(bam: os.PathLike, annotations: os.PathLike, custom_filtering: os.PathLike = None, output_prefix: os.PathLike = 'reporters', stats_prefix: os.PathLike = "", method: str = 'capture', sample_name: str = "", read_type: str = "", gzip: bool = False)
```

Removes unwanted aligned slices and identifies reporters.

Parses a BAM file and merges this with a supplied annotation to identify unwanted slices. Filtering can be tuned for Capture-C, Tri-C and Tiled-C data to ensure optimal filtering.

Common filters include:

- Removal of unmapped slices
- Removal of excluded/blacklisted slices
- Removal of non-capture fragments
- Removal of multi-capture fragments
- Removal of non-reporter fragments
- Removal of fragments with duplicated coordinates.

For specific filtering for each of the three methods see:

- *CCSliceFilter*
- *TriCSliceFilter*
- *TiledCSliceFilter*

In addition to outputting valid reporter fragments and slices separated by capture probe, this script also provides statistics on the number of read/slices filtered at each stage, and the number of cis/trans reporters for each probe.

## Notes

Whilst the script is capable of processing any annotations in tsv format, provided that the correct columns are present. It is highly recommended that the “annotate” subcommand is used to generate this file.

Slice filtering is currently hard coded into each filtering class. This will be modified in a future update to enable custom filtering orders.

## Parameters

- **bam** (*os.PathLike*) – Input bam file to analyse
- **annotations** (*os.PathLike*) – Annotations file generated by slices-annotate
- **custom\_filtering** (*os.PathLike*) – Allows for custom filtering to be performed. A yaml file is used to supply this ordering.
- **output\_prefix** (*os.PathLike, optional*) – Output file prefix. Defaults to “reporters”.
- **stats\_prefix** (*os.PathLike, optional*) – Output stats prefix. Defaults to “”.
- **method** (*str, optional*) – Analysis method. Choose from (capture|tri|tiled). Defaults to “capture”.
- **sample\_name** (*str, optional*) – Sample being processed e.g. DOX-treated\_1. Defaults to “”.
- **read\_type** (*str, optional*) – Process combined(flashed) or non-combined reads (pe) used for statistics. Defaults to “”.
- **gzip** (*bool, optional*) – Compress output with gzip. Defaults to False.

## fastq deduplicate

Created on Fri Oct 4 13:47:20 2019 @author: asmith

```
ccanalyser.cli.fastq_deduplicate.parse(input_files: Tuple, output: os.PathLike = 'out.json',
                                         read_buffer: int = 100000.0)
```

Parses fastq file(s) into easy to deduplicate format.

This command parses one or more fastq files and generates a dictionary containing hashed read identifiers together with hashed concatenated sequences. The hash dictionary is output in json format and the identify subcommand can be used to determine which read identifiers have duplicate sequences.

### Parameters

- **input\_files** (*Tuple*) – One or more fastq files to process
- **output** (*os.PathLike, optional*) – Output for parsed read identifiers and sequences. Defaults to “out.json”.
- **read\_buffer** (*int, optional*) – Number of reads to process before outputting to file. Defaults to 1e5.

```
ccanalyser.cli.fastq_deduplicate.identify(input_files: Tuple, output: os.PathLike = 'duplicates.json')
```

Identifies fragments with duplicated sequences.

Merges the hashed dictionaries (in json format) generated by the “parse” subcommand and identifies read with exactly the same sequence (share an identical hash). Duplicated read identifiers (hashed) are output in json format. The “remove” subcommand uses this dictionary to remove duplicates from fastq files.

### Parameters

- **input\_files** (*Tuple*) – Paths to json files containing dictionaries with hashed read ids as the keys and hashed sequences as the values.
- **output** (*os.PathLike, optional*) – Duplicate read ids identified. Defaults to “duplicates.json”.

```
ccanalyser.cli.fastq_deduplicate.remove(input_files: Tuple, duplicated_ids: os.PathLike,
                                         read_buffer: int = 100000.0, output_prefix: os.PathLike = '',
                                         gzip: bool = False, compression_level: int = 5, sample_name: str = '',
                                         stats_prefix: os.PathLike = '')
```

Removes fragments with duplicated sequences from fastq files.

Parses input fastq files and removes any duplicates from the fastq file(s) that are present in the json file supplied. This json dictionary should be produced by the “identify” subcommand.

Statistics for the number of duplicated and unique reads are also provided.

### Parameters

- **input\_files** (*Tuple*) – Input fastq files (in the same order as used for the parse command).
- **duplicated\_ids** (*os.PathLike*) – Duplicated read ids from identify command (hashed and in json format).
- **read\_buffer** (*int, optional*) – Number of reads to process before writing to file. Defaults to 1e5.
- **output\_prefix** (*os.PathLike, optional*) – Deduplicated fastq output prefix. Defaults to “”.

- **gzip** (*bool, optional*) – Determines if output is gzip compressed using pigz. Defaults to False.
- **compression\_level** (*int, optional*) – Level of compression if required (1-9). Defaults to 5.
- **sample\_name** (*str, optional*) – Name of sample processed e.g. DOX-treated\_1. Defaults to “”.
- **stats\_prefix** (*os.PathLike, optional*) – Output prefix for statistics. Defaults to “”.

## fastq digest

```
ccanalyser.cli.fastq_digest.collate_statistics(statq: multiprocessing.context.BaseContext.Queue,  
n_subprocesses: int) → pandas.core.frame.DataFrame
```

Collates digestion statistics from supplied statistics queue.

### Parameters

- **statq** (*Queue*) – Queue to use for collating statistics. Final item(s) must be ‘END’
- **n\_subprocesses** (*int*) – Number of digestion processes used. Required to know when to stop aquiring data from the queue.

### Returns

**Digestion statistics in histogram format.** Columns: ‘read\_type’, ‘read\_number’, ‘unfiltered/filtered’, ‘n\_slices’, ‘n\_reads’

### Return type

```
ccanalyser.cli.fastq_digest.digest(input_fastq: Tuple, restriction_enzyme: str, mode: str  
= 'pe', output_file: os.PathLike = 'out.fastq.gz', minimum_slice_length: int = 18, compression_level: int = 5,  
n_cores: int = 1, read_buffer: int = 100000, stats_prefix:  
os.PathLike = "", keep_cutsite: bool = False, sample_name:  
str = "")
```

Performs in silico digestion of one or a pair of fastq files.

### Parameters

- **input\_fastq** (*Tuple*) – Input fastq files to process
- **restriction\_enzyme** (*str*) – Restriction enzyme name or site to use for digestion.
- **mode** (*str, optional*) – Digest combined(flashed) or non-combined(pe). Undigested pe reads are output but flashed are not written. Defaults to “pe”.
- **output\_file** (*os.PathLike, optional*) – Output fastq file path. Defaults to “out.fastq.gz”.
- **minimum\_slice\_length** (*int, optional*) – Minimum allowed length for in silico digested reads. Defaults to 18.
- **compression\_level** (*int, optional*) – Compression level for gzip output (1-9). Defaults to 5.
- **n\_cores** (*int, optional*) – Number of digestion processes to use. Defaults to 1.

- **read\_buffer** (*int, optional*) – Number of reads to process before writing to file. Defaults to 100000.
- **stats\_prefix** (*os.PathLike, optional*) – Output prefix for stats file. Defaults to “”.
- **keep\_cutsite** (*bool, optional*) – Determines if cutsite is removed from the output. Defaults to False.
- **sample\_name** (*str, optional*) – Name of sample processed eg. DOX-treated\_1. Defaults to “”.

## fastq split

Created on Wed Jan 8 15:45:09 2020

@author: asmith

Script splits a fastq into specified chunks

```
ccanalyser.cli.fastq_split.run_unix_split(fn: os.PathLike, n_reads: int, read_number: int,
                                         output_prefix: os.PathLike = "", gzip: bool =
                                         False, compression_level: int = 5)

ccanalyser.cli.fastq_split.split(input_files: Tuple, method: str = 'unix', output_prefix:
                                    os.PathLike = 'split', compression_level: int = 5, n_reads: int
                                    = 1000000, gzip: bool = True)
```

Splits fastq file(s) into equal chunks of n reads.

### Parameters

- **input\_files** (*Tuple*) – Input fastq files to process.
- **method** (*str, optional*) – Python or unix method (faster but not guaranteed to maintain read pairings) to split the fastq files. Defaults to “unix”.
- **output\_prefix** (*os.PathLike, optional*) – Output prefix for split fastq files. Defaults to “split”.
- **compression\_level** (*int, optional*) – Compression level for gzipped output. Defaults to 5.
- **n\_reads** (*int, optional*) – Number of reads to split the input fastq files into. Defaults to 1000000.
- **gzip** (*bool, optional*) – Gzip compress output files if True. Defaults to True.

## genome digest

Created on Fri Oct 4 13:47:20 2019 @author: asmith

Script generates a bed file of restriction fragment locations in a given genome.

```
ccanalyser.cli.genome_digest.parse_chromosomes(fasta: pysam.libcfaidx.FastxFile) → It-
erator[pysam.libcfaidx.FastqProxy]
```

Parses a whole genome fasta file and yields chromosome entries.

**Parameters** **fasta** (*pysam.FastxFile*) – Fasta file to process.

**Yields** *Iterator[pysam.FastqProxy]* – Chromosome entry.

```
ccanalyser.cli.genome_digest.digest(input_fasta: os.PathLike, recognition_site: str, log-file: os.PathLike = 'genome_digest.log', output_file: os.PathLike = 'genome_digest.bed', remove_cutsite: bool = True, sort=False)
```

Performs in silico digestion of a genome in fasta format.

Digests the supplied genome fasta file and generates a bed file containing the locations of all restriction fragments produced by the supplied restriction enzyme.

A log file recording the number of restriction fragments for the supplied genome is also generated.

#### Parameters

- **input\_fasta** (*os.PathLike*) – Path to fasta file containing whole genome sequence, split by chromosome
- **recognition\_site** (*str*) – Restriction enzyme name/ Sequence of recognition site.
- **logfile** (*os.PathLike, optional*) – Output path of the digestion logfile. Defaults to genome\_digest.log.
- **output\_file** (*os.PathLike, optional*) – Output path for digested chromosome bed file. Defaults to genome\_digest.bed.
- **remove\_cutsite** (*bool, optional*) – Determines if restriction site is removed. Defaults to True.

### reporters count

```
ccanalyser.cli.reporters_count.count_re_site_combinations(groups: pandas.core.groupby.groupby.GroupBy, column: str = 'restriction_fragment', counts: Optional[collections.defaultdict] = None) → collections.defaultdict
```

Counts the number of unique combinations between groups in a column.

#### Parameters

- **df** (*pd.core.groupby.DataFrame*) – Aggregated dataframe for processing.
- **column** (*str, optional*) – Column to examine for unique combinations per group. Defaults to “restriction\_fragment”.
- **counts** (*defaultdict, optional*) – defaultdict(int) containing previous counts. Defaults to None.

**Returns** defaultdict(int) containing the count of unique interactions.

#### Return type

```
ccanalyser.cli.reporters_count.count(reporters: os.PathLike, output: os.PathLike = 'counts.tsv', remove_exclusions: bool = False, remove_capture: bool = False, subsample: int = 0)
```

Determines the number of captured restriction fragment interactions genome wide.

Parses a reporter slices tsv and counts the number of unique restriction fragment interaction combinations that occur within each fragment.

Options to ignore unwanted counts e.g. excluded regions or capture fragments are provided. In addition the number of reporter fragments can be subsampled if required.

### Parameters

- **reporters** (*os.PathLike*) – Reporter tsv file path.
- **output** (*os.PathLike, optional*) – Output file path for interaction counts tsv. Defaults to ‘counts.tsv’.
- **remove\_exclusions** (*bool, optional*) – Removes regions marked as capture proximity exclusions before counting. Defaults to False.
- **remove\_capture** (*bool, optional*) – Removes all capture fragments before counting. Defaults to False.
- **subsample** (*int, optional*) – Subsamples the fragments by the specified fraction. Defaults to 0 i.e. No subsampling.

### **reporters differential**

```
ccanalyser.cli.reporters_differential.get_chromosome_from_name(df:      pan-
                                                               das.core.frame.DataFrame,
                                                               name: str)

ccanalyser.cli.reporters_differential.differential(union_bedgraph: os.PathLike,
                                                   capture_name: str, cap-
                                                   ture_oligos: os.PathLike,
                                                   output_prefix: os.PathLike = 'differential', design_matrix:
                                                   Optional[os.PathLike] = None,
                                                   grouping_col: str = 'condi-
                                                   tion', threshold_count: float = 20, threshold_q: float = 0.05,
                                                   threshold_mean: float = 0)
```

Identifies differential interactions between conditions.

Parses a union bedgraph containg reporter counts from at least two conditions with two or more replicates for a single capture probe and outputs differential interaction results. Following filtering to ensure that the number of interactions is above the required threshold (–threshold\_count), diffxpy is used to run a wald test after fitting a negative binomial model to the interaction counts. The options to filter results can be filtered by a minimum mean value (threshold\_mean) and/or maximum q-value (threshold-q) are also provided.

### Notes

Currently both the capture oligos and the name of the probe being analysed must be provided in order to correctly extract cis interactions.

If a N\_SAMPLE \* METADATA design matrix has not been supplied, the script assumes that the standard replicate naming structure has been followed i.e. SAMPLE\_CONDITION\_REPLICATE\_(1|2).fastq.gz.

### Parameters

- **union\_bedgraph** (*os.PathLike*) – Union bedgraph containg all samples to be compared.
- **capture\_name** (*str*) – Name of capture probe. MUST match one probe within the supplied oligos.

- **capture\_oligos** (*os.PathLike*) – Capture oligos used for the analysis.
- **output\_prefix** (*os.PathLike, optional*) – Output prefix for differential interactions. Defaults to ‘differential’.
- **design\_matrix** (*os.PathLike, optional*) – Design matrix to use for grouping samples. (N\_SAMPLES \* METADATA). Defaults to None.
- **grouping\_col** (*str, optional*) – Column to use for grouping. Defaults to ‘condition’.
- **threshold\_count** (*float, optional*) – Minimum number of reported interactions required. Defaults to 20.
- **threshold\_q** (*float, optional*) – Maximum q-value for output. Defaults to 0.05.
- **threshold\_mean** (*float, optional*) – Minimum mean value for output. Defaults to 0.

## reporters heatmap

```
ccanalyser.cli.reporters_heatmap.plot_matrix(matrix, figsize=(10, 10), axis_labels=None, cmap=None, vmin=0, vmax=0)
```

```
ccanalyser.cli.reporters_heatmap.plot(cooler_fn: os.PathLike, coordinates: Union[str, os.PathLike], resolution: int, capture_names: Optional[Tuple] = None, normalisation: Optional[str] = None, cmap: str = 'jet', vmax: float = 1, vmin: float = 0, output_prefix: os.PathLike = "", remove_capture: bool = False)
```

Plots a heatmap of reporter interactions.

Parses a HDF5 file containing the result of a capture experiment (binned into even genomic windows) and plots a heatmap of interactions over a specified genomic range. If a capture probe name is not supplied the script will plot all probes present in the file.

**Heatmaps can also be normalised (–normalise) using either:**

- raw: No normalisation is performed.
- n\_interactions: The number of cis interactions.
- **n\_rf\_n\_interactions:** Normalised to the number of restriction fragments making up both genomic bins and by the number of cis interactions.
- **ice:** ICE normalisation followed by number of cis interactions correction.

## Parameters

- **cooler\_fn** (*os.PathLike*) – Path to capture cooler file containing interactions
- **coordinates** (*Union[str, os.PathLike]*) – Coordinates for plotting. Either chrX:1000-2000 or bed file. If a bed file, the capture probe name must be contained in the name.
- **resolution** (*int*) – Genomic resolution to plot. Must be present in the cooler.
- **capture\_names** (*Tuple, optional*) – Capture probes to plot. If None will plot all probes. Defaults to None.
- **normalisation** (*str, optional*) – Normalisation for heatmap. Choose from (n\_interactions|n\_rf\_n\_interactions|ice). Defaults to None.

- **cmap** (*str, optional*) – Colour map to use for heatmap. Defaults to “jet”.
- **vmax** (*float, optional*) – vmaxold for heatmap. Defaults to 1.
- **output\_prefix** (*os.PathLike, optional*) – Output prefix for heatmap. Defaults to “”.

## reporters pileup

```
ccanalyser.cli.reporters_pileup.bedgraph(cooler_fn: os.PathLike, capture_names: Optional[list] = None, output_prefix: os.PathLike = "", normalise: bool = False, binsize: int = 0, gzip: bool = True, scale_factor: int = 1000000.0, sparse: bool = True)
```

Extracts reporters from a capture experiment and generates a bedgraph file.

Identifies reporters for a single probe (if a probe name is supplied) or all capture probes present in a capture experiment HDF5 file.

The bedgraph generated can be normalised by the number of cis interactions for inter experiment comparisons and/or binned into even genomic windows.

### Parameters

- **cooler\_fn** (*os.PathLike*) – Path to hdf5 file containing cooler groups.
- **capture\_names** (*list, optional*) – Name of capture probe to extract. If None, will process all probes present in the file. Defaults to None.
- **output\_prefix** (*os.PathLike, optional*) – Output file prefix for bedgraph. Defaults to “”.
- **normalise** (*bool, optional*) – Normalise counts using the number of cis interactions. Defaults to False.
- **binsize** (*int, optional*) – Genomic binsize to use for generating bedgraph. No binning performed if less than 0. Defaults to 0.
- **gzip** (*bool, optional*) – Compress output bedgraph with gzip. Defaults to True.
- **scale\_factor** (*int, optional*) – Scaling factor for normalisation. Defaults to 1e6.
- **sparse** (*bool, optional*) – Produce bedgraph containing just positive bins (True) or all bins (False). Defaults to True.

## reporters store

```
ccanalyser.cli.reporters_store.fragments(counts: os.PathLike, fragment_map: os.PathLike, output: os.PathLike, capture_name: str, capture_viewpoints: os.PathLike, genome: str = "", suffix: str = "")
```

Stores restriction fragment interaction combinations at the restriction fragment level.

Parses reporter restriction fragment interaction counts produced by “ccanalyser reporters count” and gerates a cooler formatted group in an HDF5 File. See <https://cooler.readthedocs.io/en/latest/> for further details.

### Parameters

- **counts** (*os.PathLike*) – Path to restriction fragment interactions counts .tsv file.

- **fragment\_map** (*os.PathLike*) – Path to restriction fragment .bed file, generated with genome-digest command.
- **output** (*os.PathLike*) – Output file path for cooler hdf5 file.
- **capture\_name** (*str*) – Name of capture probe.
- **capture\_viewpoints** (*os.PathLike*) – Path to capture viewpoints bed file.
- **genome** (*str, optional*) – Name of genome used for alignment e.g. hg19. Defaults to “”.
- **suffix** (*str, optional*) – Suffix to append to filename. Defaults to “”.

```
ccanalyser.cli.reporters_store.bins(cooler_fn: os.PathLike, output: os.PathLike, binsizes:  
                                      Optional[Tuple] = None, normalise: bool = False,  
                                      n_cores: int = 1, scale_factor: int = 1000000.0, overlap_fraction:  
                                      float = 1e-09, conversion_tables: Optional[os.PathLike] = None)
```

Convert a cooler group containing restriction fragments to constant genomic windows

Parses a cooler group and aggregates restriction fragment interaction counts into genomic bins of a specified size. If the normalise option is selected, columns containing normalised counts are added to the pixels table of the output

## Notes

To avoid repeatedly calculating restriction fragment to bin conversions, bin conversion tables (a .pkl file containing a dictionary of *:class:ccanalyser.tools.storage.GenomicBinner* objects, one per binsize) can be supplied.

### Parameters

- **cooler\_fn** (*os.PathLike*) – Path to cooler file. Nested coolers can be specified by STORE\_FN.hdf5::/PATH\_TO\_COOLER
- **output** (*os.PathLike*) – Path for output binned cooler file.
- **binsizes** (*Tuple, optional*) – Genomic window sizes to use for binning. Defaults to None.
- **normalise** (*bool, optional*) – Normalise the number of interactions to total number of cis interactions (True). Defaults to False.
- **n\_cores** (*int, optional*) – Number of cores to use for binning. Performed in parallel by chromosome. Defaults to 1.
- **scale\_factor** (*int, optional*) – Scaling factor to use for normalising interactions. Defaults to 1e6.
- **overlap\_fraction** (*float, optional*) – Minimum fraction to use for defining overlapping bins. Defaults to 1e-9.

```
ccanalyser.cli.reporters_store.merge(coolers: Tuple, output: os.PathLike)
```

Merges ccanalyser cooler files together.

Produces a unified cooler with both restriction fragment and genomic bins whilst reducing the storage space required by hard linking the “bins” tables to prevent duplication.

### Parameters

- **coolers** (*Tuple*) – Cooler files produced by either the fragments or bins subcommands.
- **output** (*os.PathLike*) – Path from merged cooler file.

## tsv aggregate

Created on Wed Dec 11 21:49:19 2019

@author: davids, asmith

Join any number of tab delimited text files on a single column. Index column must have the same header name in each file to be joined. Performs an outer join on the index column. Assumes files have headers.

```
ccanalyser.cli.tsv_aggregate.format_index_var(var)
ccanalyser.cli.tsv_aggregate.replace_na(df)
ccanalyser.cli.tsv_aggregate.is_compressed(files)
ccanalyser.cli.tsv_aggregate.load_tsv(tsv, index=None, header=None)
ccanalyser.cli.tsv_aggregate.join_tsvs(fnames, index_col, n_processes=8, header=True)
ccanalyser.cli.tsv_aggregate.concat_tsvs(fnames, delayed=False, header=None)
ccanalyser.cli.tsv_aggregate.main(input_files, output, index=None, header=None,
                                method=None, n_processes=8, groupby_columns=None,
                                aggregate_method=None, aggregate_columns=None)
```

### 5.1.2 Pipeline

This pipeline processes data from Capture-C/NG Capture-C/Tri-C and Tiled-C sequencing protocols designed to identify 3D interactions in the genome from a specified viewpoint.

It takes Illumina paired-end sequencing reads in fastq format (gzip compression is preferred) as input and performs the following steps:

1. Identifies all restriction fragments in the genome
2. Quality control of raw reads (fastqc, multiqc)
3. Splits fastqs into smaller files to enable fast parallel processing.
4. Removal of PCR duplicates based on exact sequence matches from fastq files
5. Trimming of reads to remove adaptor sequence (trim\_galore)
6. Combining overlapping read pairs (FLASh)
7. In silico digestion of reads in fastq files
8. Alignment of fastq files with a user specified aligner (i.e. bowtie/bowtie2; BWA is not supported)
9. Analysis of alignment statistics (picard CollectAlignmentSummaryMetrics, multiqc)
10. Annotation of mapped reads with overlaps of capture probes, exclusion regions, blacklist, restriction fragments
11. Removal of non-reporter slices and identification of reporters
12. Removal of PCR duplicates (exact coordinate matches)
13. Storage of reporters in *cooler format* <<https://cooler.readthedocs.io/en/latest/datamodel.html>>
14. Generation of bedgraphs/BigWigs.
15. Collation of run statistics and generation of a run report

Optional:

- Generation of a UCSC track hub for visualisation.

- Differential interaction identification.
- Generation of subtraction bedgraphs for between condition comparisons
- Plotting of heatmaps.

@authors: asmith, dsims

`ccanalyser.pipeline.pipeline.modify_pipeline_params_dict()`

Modifies P.PARAMS dictionary.

- Selects the correct conda environment
- Ensures the correct queue manager is selected.
- Corrects the name of a UCSC hub by removing spaces and incorrect characters.

`ccanalyser.pipeline.pipeline.set_up_chromsizes()`

Ensures that genome chromsizes are present.

If chromsizes are not provided this function attempts to download them from UCSC. The P.PARAMS dictionary is updated with the location of the chromsizes.

`ccanalyser.pipeline.pipeline.genome_digest(infile, outfile)`

In silico digestion of the genome to identify restriction fragment coordinates.

Runs *ccanalyser genome digest*.

`ccanalyser.pipeline.pipeline.fastq_qc(infile, outfile)`

Runs fastqc on the input files to generate fastq statistics.

`ccanalyser.pipeline.pipeline.fastq_multiqc(infile, outfile)`

Collate fastqc reports into single report using multiqc

`ccanalyser.pipeline.pipeline.fastq_split(infiles, outfile)`

Splits the input fastq files into chunks for parallel processing

Runs *ccanalyser fastq split*.

`ccanalyser.pipeline.pipeline.fastq_duplicates_parse(infiles, outfile, sample_name, part_no)`

Parses fastq files into json format for sequence deduplication.

Runs *ccanalyser fastq deduplicate parse*

`ccanalyser.pipeline.pipeline.fastq_duplicates_identify(infiles, outfile)`

Identifies duplicate sequences from parsed fastq files in json format.

Runs *ccanalyser fastq deduplicate identify*

`ccanalyser.pipeline.pipeline.fastq_duplicates_remove(infiles, outfile)`

Removes duplicate read fragments identified from parsed fastq files.

`ccanalyser.pipeline.pipeline.stats_deduplication_collate(infiles, outfile)`

Combines deduplication statistics from fastq file partitions.

`ccanalyser.pipeline.pipeline.fastq_trim(infiles, outfile)`

Trim adaptor sequences from fastq files using trim\_galore

`ccanalyser.pipeline.pipeline.stats_trim_collate(infiles, outfile)`

Extracts and collates adapter trimming statistics from trim\_galore output

`ccanalyser.pipeline.pipeline.fastq_flash(infiles, outfile)`

Combine overlapping paired-end reads using FLASh

`ccanalyser.pipeline.pipeline.fastq_digest_combined(infile, outfile)`

In silico restriction enzyme digest of combined (flashed) read pairs

ccanalyser.pipeline.pipeline.**fastq\_digest\_non\_combined**(*infiles, outfile*)  
     In silico restriction enzyme digest of non-combined (non-flashed) read pairs

ccanalyser.pipeline.pipeline.**stats\_digestion\_collate**(*infiles, outfile*)  
     Aggregates in silico digestion statistics from fastq file partitions.

ccanalyser.pipeline.pipeline.**fastq\_preprocessing**()

ccanalyser.pipeline.pipeline.**fastq\_alignment**(*infile, outfile*)  
     Aligns in silico digested fastq files to the genome.

ccanalyser.pipeline.pipeline.**alignments\_merge**(*infiles, outfile*)  
     Combines bam files (by flashed/non-flashed status and sample).

This task simply provides an input for picard CollectAlignmentSummaryMetrics and is only used to provide overall mapping statistics. Fastq partitions are *not* combined at this stage.

ccanalyser.pipeline.pipeline.**alignments\_index**(*infile, outfile*)  
     Indexes all bam files (both partitioned and merged)

ccanalyser.pipeline.pipeline.**pre\_annotation**()

ccanalyser.pipeline.pipeline.**annotate\_make\_exclusion\_bed**(*outfile*)  
     Generates exclusion window around each capture site

ccanalyser.pipeline.pipeline.**annotate\_sort\_viewpoints**(*outfile*)  
     Sorts the capture oligos for bedtools intersect with –sorted option

ccanalyser.pipeline.pipeline.**annotate\_sort\_blacklist**(*outfile*)  
     Sorts the capture oligos for bedtools intersect with –sorted option

ccanalyser.pipeline.pipeline.**annotate\_alignments**(*infile, outfile*)  
     Annotates mapped read slices.

Slices are annotated with:  
     \* capture name  
     \* capture count  
     \* exclusion name  
     \* exclusion count  
     \* blacklist count  
     \* restriction fragment number

ccanalyser.pipeline.pipeline.**post\_annotation**()  
     Runs the pipeline until just prior to identification of reporters

ccanalyser.pipeline.pipeline.**alignments\_filter**(*infiles, outfile*)  
     Filters slices and outputs reporter slices for each capture site

ccanalyser.pipeline.pipeline.**reporters\_collate**(*infiles, outfile, \*grouping\_args*)  
     Concatenates identified reporters

ccanalyser.pipeline.pipeline.**alignments\_deduplicate\_fragments**(*infile, outfile, read\_type*)  
     Identifies duplicate fragments with the same coordinates and order.

ccanalyser.pipeline.pipeline.**alignments\_deduplicate\_slices**(*infile, outfile, sample\_name, read\_type, capture\_oligo*)  
     Removes reporters with duplicate coordinates

ccanalyser.pipeline.pipeline.**alignments\_deduplicate\_collate**(*infiles, outfile, \*grouping\_args*)  
     Final collation of reporters by sample and capture probe

ccanalyser.pipeline.pipeline.**stats\_alignment\_filtering\_collate**(*infiles, outfile*)  
     Combination of all reporter identification and filtering statistics

ccanalyser.pipeline.pipeline.**post\_ccanalyser\_analysis**()  
     Reporters have been identified, deduplicated and collated by sample/capture probe

ccanalyser.pipeline.pipeline.**reporters\_count** (*infile, outfile*)

Counts the number of interactions identified between reporter restriction fragments

ccanalyser.pipeline.pipeline.**reporters\_store\_restriction\_fragment** (*infile, outfile, sample\_name, capture\_name*)

Stores restriction fragment interaction counts in cooler format

ccanalyser.pipeline.pipeline.**generate\_bin\_conversion\_tables** (*outfile*)

Converts restriction fragments to genomic bins.

Binning restriction fragments into genomic bins takes a substantial amount of time and memory. To avoid repeatedly performing the same action, bin conversion tables are calculated once for each required resolution and then stored as a pickle file.

ccanalyser.pipeline.pipeline.**reporters\_store\_binned** (*infile, outfile, capture\_name*)

Converts a cooler file of restriction fragments to even genomic bins.

ccanalyser.pipeline.pipeline.**reporters\_store\_merged** (*infiles, outfile, sample\_name*)

Combines cooler files together

ccanalyser.pipeline.pipeline.**pipeline\_merge\_stats** (*infiles, outfile*)

Generates a summary statistics file for the pipeline run.

ccanalyser.pipeline.pipeline.**pipeline\_make\_report** (*infile, outfile*)

Run jupyter notebook for reporting and plotting pipeline statistics

ccanalyser.pipeline.pipeline.**reporters\_make\_bedgraph** (*infile, outfiles, sample\_name*)

Extract reporters in bedgraph format from stored interactions

ccanalyser.pipeline.pipeline.**reporters\_make\_bedgraph\_normalised** (*infile, outfiles, sample\_name*)

Extract reporters in bedgraph format from stored interactions.

In addition to generating a bedgraph this task also normalises the counts by the number of cis interactions identified to enable cross sample comparisons.

ccanalyser.pipeline.pipeline.**reporters\_make\_union\_bedgraph** (*infiles, outfile, normalisation\_type, capture\_name*)

Collates bedgraphs by capture probe into a single file for comparison.

See [bedtools unionbedg](#) for more details.

ccanalyser.pipeline.pipeline.**reporters\_make\_subtraction\_bedgraph** (*infile, outfile*)

Uses UCSC tools bedGraphToBigWig to generate bigWigs for each bedgraph

ccanalyser.pipeline.pipeline.**hub\_make** (*infiles, outfile, statistics*)

Creates a ucsc hub from the pipeline output

ccanalyser.pipeline.pipeline.**hub\_write\_path** (*outfile*)

Convinence task to write hub url to use for adding custom hub to UCSC genome browser

ccanalyser.pipeline.pipeline.**identify\_differential\_interactions** (*infile, outfile, capture\_name*)

---

```
ccanalyser.pipeline.pipeline.reporters_plot_heatmap(infile, outfile)
    Plots a heatmap over a specified region
ccanalyser.pipeline.pipeline.full(infiles, outfile)
```

### 5.1.3 CCanalyser tools

#### ccanalyser.tools.annotate module

```
class ccanalyser.tools.annotate.BedIntersection(bed1: Union[str, pybedtools.bedtool.BedTool, pandas.core.frame.DataFrame], bed2: Union[str, pybedtools.bedtool.BedTool, pandas.core.frame.DataFrame], intersection_name: str = 'count', intersection_method: str = 'count', intersection_min_frac: float = 1e-09, intersection_split_chrom: bool = True, n_cores: int = 1, invalid_bed_action='error')
```

Bases: object

Performs intersection between two named bed files.

Wrapper around bedtools intersect designed to intersect in parallel (by splitting file based on chromosome) and handle malformed bed files.

##### **bed1**

Bed file to intersect. Must be named.

**Type** Union[str, BedTool, pd.DataFrame]

##### **bed2**

Bed file to intersect.

**Type** Union[str, BedTool, pd.DataFrame]

##### **intersection\_name**

Name for intersection.

**Type** str

##### **min\_frac**

Minimum fraction required for intersection

**Type** float

##### **n\_cores**

Number of cores for parallel intersection.

**Type** int

##### **invalid\_bed\_action**

Method to deal with missing/malformed bed files ("ignore"|"error")

##### **property intersection**

Intersects the two bed files and returns a pd.Series.

## ccanalyser.tools.deduplicate module

```
class ccanalyser.tools.deduplicate.ReadDeduplicationParserProcess(inq: multiprocessing.Queue,
                                                               outq: multiprocessing.Queue,
                                                               hash_seed: int = 42,
                                                               save_hashed_dict_path: os.PathLike
                                                               = 'parsed.json')
```

Bases: multiprocessing.context.Process

Process subclass for parsing fastq file(s) into a hashed {id:sequence} json format.

### inq

Input read queue

### outq

Output read queue (Not currently used)

### hash\_seed

Seed for xxhash64 algorithm to ensure consistency

### save\_hashed\_dict\_path

Path to save hashed dictionary

### run()

Processes fastq reads from multiple files and generates a hashed json dictionary.

Dictionary is hashed and in the format {(read 1 name + read 2 name): (sequence 1 + sequence 2)}.

Output path is specified by save\_hashed\_dict\_path.

```
class ccanalyser.tools.deduplicate.ReadDuplicateRemovalProcess(inq: multiprocessing.Queue,
                                                               outq: multiprocessing.Queue,
                                                               duplicated_ids: set,
                                                               statq: Optional[multiprocessing.context.BaseContext.SimpleQueue],
                                                               = None,
                                                               hash_seed: int = 42)
```

Bases: multiprocessing.context.Process

Process subclass for parsing fastq file(s) and removing identified duplicates.

### inq

Input read queue

### outq

Output queue for deduplicated reads.

### duplicated\_ids

Concatenated read ids to remove from input fastq files.

---

**statq**  
Output queue for statistics.

**reads\_total**  
Number of fastq reads processed.

**reads\_unique**  
Number of non-duplicated reads output.

**hash\_seed**  
Seed for xxhash algorithm. MUST be the same as used by ReadDuplicationParserProcess.

**run()**  
Performs read deduplication based on sequence.  
Unique reads are placed on outq and deduplication statistics are placed on statq.

## ccanalyser.tools.digest module

**class** ccanalyser.tools.digest.DigestedChrom(*chrom*: pysam.libcfaidx.FastqProxy, *cutsite*: str, *fragment\_number\_offset*: int = 0, *fragment\_min\_len*: int = 1)

Bases: object

Performs in silico digestion of fasta files.

Identifies all restriction sites for a supplied restriction enzyme/restriction site and generates bed style entries.

**chrom**  
Chromosome to digest  
**Type** pysam.FastqProxy

**recognition\_seq**  
Sequence of restriction recognition site  
**Type** str

**recognition\_len**  
Length of restriction recognition site  
**Type** int

**recognition\_seq**  
Regular expression for restriction recognition site  
**Type** re.Pattern

**fragment\_indexes**  
Indexes of fragment(s) start and end positions.  
**Type** List[int]

**fragment\_number\_offset**  
Starting fragment number.  
**Type** int

**fragment\_min\_len**  
Minimum fragment length required to report fragment  
**Type** int

**get\_recognition\_site\_indexes()** → List[int]  
Gets the start position of all recognition sites present in the sequence.

## Notes

Also appends the start and end of the sequence to enable clearer iteration through the indexes.

**Returns** Indexes of fragment(s) start and end positions.

**Return type** List[int]

### property fragments

Extracts the coordinates of restriction fragments from the sequence.

**Yields** Iterator[Iterable[str]] – Identified restriction fragments in bed format.

```
class ccanalyser.tools.digest.DigestedRead(read: pysam.libcfaidx.FastqProxy, cut-site: str, min_slice_length: int = 18, slice_number_offset: int = 0, allow_undigested: bool = False, read_type: str = 'flashed')
```

Bases: object

Performs in slice digestion of fastq files.

Identifies all restriction sites for a supplied restriction enzyme/restriction site and generates bed style entries.

### read

Read to digest.

**Type** pysam.FastqProxy

### recognition\_seq

Sequence of restriction recognition site.

**Type** str

### recognition\_len

Length of restriction recognition site.

**Type** int

### recognition\_seq

Regular expression for restriction recognition site.

**Type** re.Pattern

### slices

List of Fastq formatted digested reads (slices).

**Type** List[str]

### slice\_indexes

Indexes of fragment(s) start and end positions.

**Type** List[int]

### slice\_number\_offset

Starting fragment number.

**Type** int

### min\_slice\_len

Minimum fragment length required to report fragment.

**Type** int

### has\_slices

Recognition site(s) present within sequence.

**Type** bool

**get\_recognition\_site\_indexes()** → List[int]

---

**class** ccanalyser.tools.digest.**ReadDigestionProcess** (*inq: multiprocessing.context.BaseContext.SimpleQueue, outq: multiprocessing.context.BaseContext.SimpleQueue, statq: Optional[multiprocessing.context.BaseContext.Queue] = None, \*\*digestion\_kwargs*)

Bases: multiprocessing.context.Process

Process subclass for multiprocessing fastq digestion.

**run()**

Performs read digestion.

Reads to digest are pulled from inq, digested with the DigestedRead class and the results placed on outq for writing.

If a statq is provided, read digestion stats are placed into this queue for aggregation.

## ccanalyser.tools.filter module

---

**class** ccanalyser.tools.filter.**SliceFilter** (*slices: pandas.core.frame.DataFrame, filter\_stages: Optional[dict] = None, sample\_name: str = "", read\_type: str = ""*)

Bases: object

Perform slice filtering (inplace) and reporter identification.

The SliceFilter classes e.g. CCSliceFilter, TriCSliceFilter, TiledCSliceFilter perform all of the filtering (inplace) and reporter identification whilst also providing statistics of the numbers of slices/reads removed at each stage.

### **slices**

Type pd.DataFrame

### **fragments**

Slices dataframe aggregated by parental read.

Type pd.DataFrame

### **reporters**

Slices identified as reporters.

Type pd.DataFrame

### **filter\_stages**

Dictionary containing stages and a list of class methods (str) required to get to this stage.

Type dict

### **slice\_stats**

Provides slice level statistics.

Type pd.DataFrame

### **read\_stats**

Provides statistics of slice filtering at the parental read level.

Type pd.DataFrame

**filter\_stats**

Provides statistics of read filtering.

**Type** pd.DataFrame

**property filters**

A list of the callable filters present within the slice filterer instance.

**Returns** All filters present in the class.

**Return type** list

**property slice\_stats**

Statistics at the slice level.

**Returns** Statistics per slice.

**Return type** pd.DataFrame

**property filter\_stats**

Statistics for each filter stage.

**Returns** Statistics of the number of slices removed at each stage.

**Return type** pd.DataFrame

**property read\_stats**

Gets statistics at a read level.

Aggregates slices by parental read id and calculates stats.

**Returns** Statistics of the slices/fragments removed aggregated by read id.

**Return type** pd.DataFrame

**property fragments**

Summarises slices at the fragment level.

Uses pandas groupby to aggregate slices by their parental read name (shared by all slices from the same fragment). Also determines the number of reporter slices for each fragment.

**Returns** Slices aggregated by parental read name.

**Return type** pd.DataFrame

**property reporters**

Extracts reporter slices from slices dataframe i.e. non-capture slices

**Returns** All non-capture slices

**Return type** pd.DataFrame

**filter\_slices (output\_slices=False, output\_location='.)'**

Performs slice filtering.

Filters are applied to the slices dataframe in the order specified by filter\_stages. Filtering stats aggregated at the slice and fragment level are also printed.

**Parameters**

- **output\_slices (bool, optional)** – Determines if slices are to be output to a specified location after each filtering step. Useful for debugging. Defaults to False.
- **output\_location (str, optional)** – Location to output slices at each stage. Defaults to “.”.

---

**get\_unfiltered\_slices()**  
Does not modify slices.

**remove\_unmapped\_slices()**  
Removes slices marked as unmapped (Uncommon)

**remove\_orphan\_slices()**  
Remove fragments with only one aligned slice (Common)

**remove\_duplicate\_re frags()**  
Prevent the same restriction fragment being counted more than once (Uncommon).

### Example

-RE\_FRAG1---Capture----RE\_FRAG1--

**remove\_slices\_without\_re\_frag\_assigned()**  
Removes slices if restriction\_fragment column is N/A

**remove\_duplicate\_slices()**  
Remove all slices if the slice coordinates and slice order are shared.  
This method is designed to remove a fragment if it is a PCR duplicate (Common).

### Example

Frag 1: chr1:1000-1250 chr1:1500-1750  
 Frag 2: chr1:1000-1250 chr1:1500-1750  
 Frag 3: chr1:1050-1275 chr1:1600-1755  
 Frag 4: chr1:1500-1750 chr1:1000-1250

Frag 2 removed. Frag 1,3,4 retained

**remove\_duplicate\_slices\_pe()**  
Removes PCR duplicates from non-flashed (PE) fragments (Common).  
Sequence quality is often lower at the 3' end of reads leading to variance in mapping coordinates. PCR duplicates are removed by checking that the fragment start and end are not duplicated in the dataframe.

**remove\_excluded\_slices()**  
Removes any slices in the exclusion region (default 1kb) (V. Common)

**remove\_blacklisted\_slices()**  
Removes slices marked as being within blacklisted regions

**class ccanalyser.tools.filter.CCSliceFilter(slices, filter\_stages=None, \*\*sample\_kwargs)**  
Bases: *ccanalyser.tools.filter.SliceFilter*

Perform Capture-C slice filtering (inplace) and reporter identification.

SliceFilter tuned specifically for Capture-C data. This class has additional methods to remove common artifacts in Capture-C data i.e. multi-capture fragments, non-reporter fragments, multi-capture reporters. The default filter order is as follows:

- remove\_unmapped\_slices
- remove\_orphan\_slices

- remove\_multi\_capture\_fragments
- remove\_excluded\_slices
- remove\_blacklisted\_slices
- remove\_non\_reporter\_fragments
- remove\_multicapture\_reporters
- remove\_slices\_without\_re\_frag\_assigned
- remove\_duplicate\_re frags
- remove\_duplicate\_slices
- remove\_duplicate\_slices\_pe
- remove\_non\_reporter\_fragments

See the individual methods for further details.

**slices**

Annotated slices dataframe.

**Type** pd.DataFrame

**fragments**

Slices dataframe aggregated by parental read.

**Type** pd.DataFrame

**reporters**

Slices identified as reporters.

**Type** pd.DataFrame

**filter\_stages**

Dictionary containing stages and a list of class methods (str) required to get to this stage.

**Type** dict

**slice\_stats**

Provides slice level statistics.

**Type** pd.DataFrame

**read\_stats**

Provides statistics of slice filtering at the parental read level.

**Type** pd.DataFrame

**filter\_stats**

Provides statistics of read filtering.

**Type** pd.DataFrame

**property\_fragments**

Summarises slices at the fragment level.

Uses pandas groupby to aggregate slices by their parental read name (shared by all slices from the same fragment). Also determines the number of reporter slices for each fragment.

**Returns** Slices aggregated by parental read name.

**Return type** pd.DataFrame

**property slice\_stats**

Statistics at the slice level.

**Returns** Statistics per slice.

**Return type** pd.DataFrame

**property frag\_stats**

Statistics aggregated at the fragment level.

As this involves slice aggregation it can be rather slow for large datasets. It is recommended to only use this property if it is required.

**Returns** Fragment level statistics

**Return type** pd.DataFrame

**property reporters**

Extracts reporter slices from slices dataframe i.e. non-capture slices

**Returns** All non-capture slices

**Return type** pd.DataFrame

**property captures**

Extracts capture slices from slices dataframe

i.e. slices that do not have a null capture name

**Returns** Capture slices

**Return type** pd.DataFrame

**property capture\_site\_stats**

Extracts the number of unique capture sites.

**property merged\_captures\_and\_reporters**

Merges captures and reporters sharing the same parental id.

Capture slices and reporter slices with the same parental read id are merged together. The prefixes ‘capture’ and ‘reporter’ are used to identify slices marked as either captures or reporters.

**Returns** Merged capture and reporter slices

**Return type** pd.DataFrame

**property cis\_or\_trans\_stats**

Extracts reporter cis/trans statistics from slices.

**Returns** Reporter cis/trans statistics

**Return type** pd.DataFrame

**remove\_non\_reporter\_fragments()**

Removes the fragment if it has no reporter slices present (Common)

**remove\_multi\_capture\_fragments()**

Removes double capture fragments.

All slices (i.e. the entire fragment) are removed if more than one capture probe is present i.e. a double capture (V. Common)

**remove\_multicapture\_reporters (n\_adjacent: int = 1)**

Deals with an odd situation in which a reporter spanning two adjacent capture sites is not removed.

## Example

—Capture 1—/—Capture 2— —REP—

In this case the “reporter” slice is not considered either a capture or exclusion.

These cases are dealt with by explicitly removing reporters on restriction fragments adjacent to capture sites.

**Parameters** `n_adjacent` – Number of adjacent restriction fragments to remove

```
class ccanalyser.tools.filter.TriCSliceFilter(slices,    filter_stages=None,      **sam-
                                                ple_kwargs)
Bases: ccanalyser.tools.filter.CCSliceFilter
```

Perform Tri-C slice filtering (inplace) and reporter identification.

SliceFilter tuned specifically for Tri-C data. Whilst the vast majority of filters are inherited from CCSliceFilter, this class has additional methods for Tri-C analysis i.e. `remove_slices_with_one_reporter`. The default filtering order is:

- `remove_unmapped_slices`
- `remove_slices_without_re_frag_assigned`
- `remove_orphan_slices`
- `remove_multi_capture_fragments`
- `remove_blacklisted_slices`
- `remove_non_reporter_fragments`
- `remove_multicapture_reporters`
- `remove_duplicate_re_frags`
- `remove_duplicate_slices`
- `remove_duplicate_slices_pe`
- `remove_non_reporter_fragments`
- `remove_slices_with_one_reporter`

See the individual methods for further details.

### `slices`

Annotated slices dataframe.

**Type** pd.DataFrame

### `fragments`

Slices dataframe aggregated by parental read.

**Type** pd.DataFrame

### `reporters`

Slices identified as reporters.

**Type** pd.DataFrame

### `filter_stages`

Dictionary containing stages and a list of class methods (str) required to get to this stage.

**Type** dict

**slice\_stats**  
Provides slice level statistics.

**Type** pd.DataFrame

**read\_stats**  
Provides statistics of slice filtering at the parental read level.

**Type** pd.DataFrame

**filter\_stats**  
Provides statistics of read filtering.

**Type** pd.DataFrame

**remove\_slices\_with\_one\_reporter()**  
Removes fragments if they do not contain at least two reporters.

**class** ccanalyser.tools.filter.TiledCSliceFilter(slices, filter\_stages=None, \*\*sample\_kwargs)  
Bases: *ccanalyser.tools.filter.SliceFilter*

Perform Tiled-C slice filtering (inplace) and reporter identification.

SliceFilter tuned specifically for Tiled-C data. This class has additional methods to remove common artifacts in Tiled-C data i.e. non-capture fragments, multi-capture (with different tiled regions) fragments. A reporter is defined differently in a Tiled-C analysis as a reporter slice can also be a capture slice.

The default filter order is as follows:

- remove\_unmapped\_slices
- remove\_orphan\_slices
- remove\_blacklisted\_slices
- remove\_non\_capture\_fragments
- remove\_dual\_capture\_fragments
- remove\_slices\_without\_re\_frag\_assigned
- remove\_duplicate\_re frags
- remove\_duplicate\_slices
- remove\_duplicate\_slices\_pe
- remove\_orphan\_slices

See the individual methods for further details.

**slices**  
Annotated slices dataframe.

**Type** pd.DataFrame

**fragments**  
Slices dataframe aggregated by parental read.

**Type** pd.DataFrame

**reporters**  
Slices identified as reporters.

**Type** pd.DataFrame

**filter\_stages**

Dictionary containing stages and a list of class methods (str) required to get to this stage.

**Type** dict

**slice\_stats**

Provides slice level statistics.

**Type** pd.DataFrame

**read\_stats**

Provides statistics of slice filtering at the parental read level.

**Type** pd.DataFrame

**filter\_stats**

Provides statistics of read filtering.

**Type** pd.DataFrame

**property fragments**

Summarises slices at the fragment level.

Uses pandas groupby to aggregate slices by their parental read name (shared by all slices from the same fragment). Also determines the number of reporter slices for each fragment.

**Returns** Slices aggregated by parental read name.

**Return type** pd.DataFrame

**property slice\_stats**

Statistics at the slice level.

**Returns** Statistics per slice.

**Return type** pd.DataFrame

**property cis\_or\_trans\_stats**

Extracts reporter cis/trans statistics from slices.

Unlike Capture-C/Tri-C reporter slice can also be capture slices as all slices within the capture region are considered as reporters. To extract cis/trans statistics, one capture slice in each fragment is considered to be the “primary capture” this then enables merging of this “primary capture” with the other reporters both inside and outside of the tiled region.

**Returns** Reporter cis/trans statistics

**Return type** pd.DataFrame

**remove\_slices\_outside\_capture()**

Removes slices outside of capture region(s)

**remove\_non\_capture\_fragments()**

Removes fragments without a capture assigned

**remove\_dual\_capture\_fragments()**

Removes a fragment with multiple different capture sites.

Modified for TiledC filtering as the fragment dataframe is generated slightly differently.

## ccanalyser.tools.io module

```
class ccanalyser.tools.io.FastqReaderProcess (input_files: Union[str, list], outq: multiprocessing.contextBaseContext.Queue, read_buffer: int = 100000, read_counter: Optional[multiprocessing.managers.BaseManager.register.<locals>.ter = None, n_subprocesses: int = 1, statq: Optional[multiprocessing.contextBaseContext.Queue] = None)
```

Bases: multiprocessing.context.Process

Reads fastq file(s) in chunks and places them on a queue.

**input\_file**  
Input fastq files.

**outq**  
Output queue for chunked reads/read pairs.

**statq**  
(Not currently used) Queue for read statistics if required.

**read\_buffer**  
Number of reads to process before placing them on outq

**read\_counter**  
(Not currently used) Can be used to sync between multiple readers.

**n\_subprocesses**  
Number of processes running concurrently. Used to make sure enough termination signals are used.

**run()**  
Performs reading and chunking of fastq file(s).

```
class ccanalyser.tools.io.FastqReadFormatterProcess (inq: multiprocess-  
ing.contextBaseContext.SimpleQueue, outq: multiprocess-  
ing.contextBaseContext.SimpleQueue, formatting: Optional[list] = None)
```

Bases: multiprocessing.context.Process

**run()**  
Method to be run in sub-process; can be overridden in sub-class

```
class ccanalyser.tools.io.FastqWriterSplitterProcess (inq: multiprocess-  
ing.contextBaseContext.Queue, output_prefix: Union[str, list], paired_output: bool = False, gzip=False, compression_level: int = 3, compression_threads: int = 8, n_subprocesses: int = 1, n_workers_terminated: int = 0, n_files_written: int = 0)
```

Bases: multiprocessing.context.Process

**run()**  
Method to be run in sub-process; can be overridden in sub-class

```
class ccanalyser.tools.io.FastqWriterProcess (inq: multiprocessing.context.BaseContext.Queue, output: Union[str, list], compression_level: int = 5, n_subprocesses: int = 1)
```

Bases: multiprocessing.context.Process

**run ()**

Method to be run in sub-process; can be overridden in sub-class

```
ccanalyser.tools.io.parse_alignment (aln)
```

Parses reads from a bam file into a list.

**Extracts:** -read name -parent reads -flashed status -slice number -mapped status -multimapping status -chromosome number (e.g. chr10) -start (e.g. 1000) -end (e.g. 2000) -coords e.g. (chr10:1000-2000)

**Parameters** **aln** – pysam.AlignmentFile.

**Returns** Containing the attributes extracted.

**Return type** list

```
ccanalyser.tools.io.parse_bam (bam)
```

Uses parse\_alignment function convert bam file to a dataframe.

**Extracts:** -‘slice\_name’ -‘parent\_read’ -‘pe’ -‘slice’ -‘mapped’ -‘multimapped’ -‘chrom’ -‘start’ -‘end’ -‘coordinates’

**Parameters** **bam** – File name of bam file to process.

**Returns** DataFrame with the columns listed above.

**Return type** pd.DataFrame

## ccanalyser.tools.pileup module

```
class ccanalyser.tools.pileup.CoolerBedGraph (cooler_fn: str, sparse: bool = True, only_cis: bool = False)
```

Bases: object

Generates a bedgraph file from a cooler file created by interactions-store.

**cooler**

Cooler file to use for bedgraph production

**Type** cooler.Cooler

**capture\_name**

Name of capture probe being processed.

**Type** str

**sparse**

Only output bins with interactions.

**Type** bool

**only\_cis**

Only output cis interactions.

**Type** bool

**property bedgraph**

Returns: pd.DataFrame: DataFrame in bedgraph format.

**property reporters**

Interactions with capture fragments/bins.

**Returns** DataFrame containing just bins interacting with the capture probe.

**Return type** pd.DataFrame

**normalise\_bedgraph (scale\_factor=1000000.0) → pandas.core.frame.DataFrame**

Normalises the bedgraph.

Uses the number of cis interactions to normalise the bedgraph counts.

**Parameters** **scale\_factor** (int, optional) – Scaling factor for normalisation. Defaults to 1e6.

**Returns** Normalised bedgraph formatted DataFrame

**Return type** pd.DataFrame

**to\_file (fn: os.PathLike, normalise: bool = False, \*\*normalise\_kwargs)**

Outputs the bedgraph dataframe to a file.

If normalise is True, will also normalise the counts by the number of cis interactions.

**Parameters**

- **fn** (os.PathLike) – Output file name.
- **normalise** (bool, optional) – Normalise the bedgraph before writing to file. Defaults to False.

```
class ccanalyser.tools.pileup.CoolerBedGraphWindowed(cooler_fn: str, binsize: int
= 5000.0, binner: Optional[ccanalyser.tools.storage.CoolerBinner]
= None, sparse=True)
```

Bases: *ccanalyser.tools.pileup.CoolerBedGraph*

**normalise\_bedgraph (scale\_factor=1000000.0)**

Normalises the bedgraph.

Uses the number of cis interactions to normalise the bedgraph counts.

**Parameters** **scale\_factor** (int, optional) – Scaling factor for normalisation. Defaults to 1e6.

**Returns** Normalised bedgraph formatted DataFrame

**Return type** pd.DataFrame

**property reporters\_binned**

```
class ccanalyser.tools.pileup.CCBedgraph(path=None, df=None, capture_name='', capture_chrom='', capture_start='', capture_end='')
```

Bases: object

**property score****property coordinates****to\_bedtool()****to\_file (path)**

**ccanalyser.tools.plotting module**

```
class ccanalyser.tools.plotting.CCMATRIX(cooler_fn: os.PathLike, binsize: 5000, cap-  
ture_name: str, remove_capture=False)  
Bases: object  
get_matrix(coordinates, field='count')  
get_matrix_normalised(coordinates, normalisation_method=None, **normalisation_kwargs)
```

**ccanalyser.tools.statistics module**

```
class ccanalyser.tools.statistics.DeduplicationStatistics(sample: str, read_type:  
str = 'pe', reads_total:  
int = 0, reads_unique:  
int = 0)  
Bases: object  
property df  
class ccanalyser.tools.statistics.DigestionStats(read_type, read_number, unfiltered,  
filtered)  
Bases: tuple  
filtered  
Alias for field number 3  
read_number  
Alias for field number 1  
read_type  
Alias for field number 0  
unfiltered  
Alias for field number 2  
ccanalyser.tools.statistics.collate_histogram_data(fnames)  
ccanalyser.tools.statistics.collate_read_data(fnames)  
ccanalyser.tools.statistics.collate_slice_data(fnames)  
ccanalyser.tools.statistics.collate_cis_trans_data(fnames)  
ccanalyser.tools.statistics.extract_trimming_stats(fn)
```

**ccanalyser.tools.storage module**

```
ccanalyser.tools.storage.get_capture_coords(viewpoint_file: str, viewpoint_name: str)  
ccanalyser.tools.storage.get_capture_bins(bins, viewpoint_chrom, viewpoint_start, view-  
point_end)  
ccanalyser.tools.storage.create_cooler_cc(output_prefix: str, bins: pandas.core.frame.DataFrame,  
pixels: pandas.core.frame.DataFrame, capture_name: str,  
capture_viewpoints: os.PathLike, capture_bins: Optional[Union[int, list]] = None, suffix=None,  
**cooler_kwargs) → os.PathLike  
Creates a cooler hdf5 file or cooler formatted group within a hdf5 file.
```

**Parameters**

- **output\_prefix** (*str*) – Output path for hdf5 file. If this already exists, will append a new group to the file.
- **bins** (*pd.DataFrame*) – DataFrame containing the genomic coordinates of all bins in the pixels table.
- **pixels** (*pd.DataFrame*) – DataFrame with columns: bin1\_id, bin2\_id, count.
- **capture\_name** (*str*) – Name of capture probe to store.
- **capture\_viewpoints** (*os.PathLike*) – Path to capture viewpoints used for the analysis.
- **capture\_bins** (*Union[int, list, optional]*) – Bins containing capture viewpoints. Can be determined from viewpoints if not supplied. Defaults to None.
- **suffix** (*str, optional*) – Suffix to append before the .hdf5 file extension. Defaults to None.

**Raises** `ValueError` – Capture name must exactly match the name of a supplied capture viewpoint.

**Returns** Path of cooler hdf5 file.

**Return type** `os.PathLike`

```
class ccAnalyser.tools.storage.GenomicBinner(chromsizes: Union[os.PathLike, pandas.core.frame.DataFrame, pandas.core.series.Series], fragments: pandas.core.frame.DataFrame, binsize: int = 5000, n_cores: int = 8, method: Literal[midpoint, overlap] = 'midpoint', min_overlap: float = 0.2)
```

Bases: `object`

Provides a conversion table for converting two sets of bins.

#### **chromsizes**

Series indexed by chromosome name containing chromosome sizes in bp

**Type** `pd.Series`

#### **fragments**

DataFrame containing bins to convert to equal genomic intervals

**Type** `pd.DataFrame`

#### **binsize**

Genomic bin size

**Type** `int`

#### **min\_overlap**

Minimum degree of intersection to define an overlap.

**Type** `float`

#### **n\_cores**

Number of cores to use for bin intersection.

**Type** `int`

#### **property bins**

Equal genomic bins.

**Returns** DataFrame in bed format.

**Return type** pd.DataFrame

**property bin\_conversion\_table**  
Returns: pd.DataFrame: Conversion table containing coordinates and ids of intersecting bins.

**class** ccanalyser.tools.storage.CoolerBinner(*cooler\_fn*: os.PathLike, *bin\_size*: Optional[int] = None, *n\_cores*: int = 8, *binner*: Optional[ccanalyser.tools.storage.GenomicBinner] = None)

Bases: object

Bins a cooler file into equal genomic intervals.

**cooler**  
(cooler.Cooler): Cooler instance to bin.

**binner**  
Binner class to generate bin conversion tables

**Type** ccanalyser.storeage.GenomicBinner

**binsize**  
Genomic bin size

**Type** int

**scale\_factor**  
Scaling factor for normalising interaction counts.

**Type** int

**n\_cis\_interactions**  
Number of cis interactions with the capture bins.

**Type** int

**n\_cores**  
Number of cores to use for binning.

**Type** int

**property bins**  
Returns: pd.DataFrame: Even genomic bins of a specified binsize.

**property bin\_conversion\_table**  
Returns: pd.DataFrame: Conversion table containing coordinates and ids of intersecting bins.

**property capture\_bins**  
Returns: pd.DataFrame: Capture bins converted to the new even genomic bin format.

**property pixel\_conversion\_table**  
Returns: pd.DataFrame: Conversion table to convert old binning scheme to the new.

**property pixels**  
Returns: pd.DataFrame: Pixels (interaction counts) converted to the new binning scheme.

**normalise\_pixels** (*n\_fragment\_correction*: bool = True, *n\_interaction\_correction*: bool = True, *scale\_factor*: int = 1000000.0)  
Normalises pixels (interactions).

Normalises pixels according to the number of restriction fragments per bin and the number of cis interactions. If both normalisation options are selected, will also provide a dual normalised column.

### Parameters

- **n\_fragment\_correction** (*bool, optional*) – Updates the pixels DataFrame with counts corrected for the number of restriction fragments per bin. Defaults to True.
- **n\_interaction\_correction** (*bool, optional*) – Updates the pixels DataFrame with counts corrected for the number of cis interactions. Defaults to True.
- **scale\_factor** (*int, optional*) – Scaling factor for n\_interaction\_correction. Defaults to 1e6.

**to\_cooler** (*store, normalise=False, \*\*normalise\_options*)

ccanalyser.tools.storage.link\_bins (*clr: os.PathLike*)

Reduces cooler storage space by linking “bins” table.

All of the cooler “bins” tables containing the genomic coordinates of each bin are identical for all cooler files of the same resolution. As cooler.create\_cooler generates a new bins table for each cooler, this leads to a high degree of duplication.

This function hard links the bins tables for a given resolution to reduce the degree of duplication.

**Parameters** **clr** (*os.PathLike*) – Path to cooler hdf5 produced by the merge command.

## 5.1.4 Utilities

ccanalyser.utils.open\_logfile (*fn: str*) → IO

Handles instances where the log file is sys.stdout

ccanalyser.utils.merge\_dictionaries (*dicts: list*) → dict

Merges multiple dictionary entries

ccanalyser.utils.invert\_dict (*d: dict*) → dict

Inverts key: value pairs into value: key pairs

ccanalyser.utils.is\_on (*param: str*) → bool

Returns True if parameter in “on” values

ccanalyser.utils.is\_off (*param: str*)

Returns True if parameter in “off” values

ccanalyser.utils.is\_none (*param: str*) → bool

Returns True if parameter is none

ccanalyser.utils.get\_human\_readable\_number\_of\_bp (*bp: int*) → str

Converts integer into human readable basepair number

ccanalyser.utils.is\_valid\_bed (*bed: Union[str, pybedtools.bedtool.BedTool], verbose=True*) → bool

Returns true if bed file can be opened and has at least 3 columns

ccanalyser.utils.bed\_has\_name (*bed: Union[str, pybedtools.bedtool.BedTool]*) → bool

Returns true if bed file has at least 4 columns

ccanalyser.utils.bed\_has\_duplicate\_names (*bed*) → bool

Returns true if bed file has no duplicated names

ccanalyser.utils.get\_re\_site (*recognition\_site: Optional[str] = None*) → str

Obtains the recognition sequence for a supplied restriction enzyme or correctly formats a supplied recognition sequence.

**Parameters**

- – **DNA sequence to use for fasta digestion e.g. "GATC"**  
(*cut\_sequence*) –
- – **Name of restriction enzyme e.g. DpnII** (*restriction\_enzyme*) –

**Returns** recognition sequence e.g. “GATC”

**Raises ValueError** – Error if restriction\_enzyme is not in known enzymes

ccanalyser.utils.**hash\_column** (*col*: Iterable, *hash\_type*=64) → list

Convinience function to perform hashing using xxhash on an iterable. Not vectorised.

ccanalyser.utils.**split\_intervals\_on\_chrom** (*intervals*: Union[str, pybedtools.bedtool.BedTool, pandas.core.frame.DataFrame]) → dict

Creates dictionary from bed file with the chroms as keys

ccanalyser.utils.**intersect\_bins** (*bins\_1*: pandas.core.frame.DataFrame, *bins\_2*: pandas.core.frame.DataFrame, \*\**bedtools\_kwargs*) → pandas.core.frame.DataFrame

Intersects two sets of genomic intervals using bedtools intersect.

Formats the intersection in a clearer way than pybedtool auto names.

ccanalyser.utils.**load\_json** (*fn*) → dict

Convinence function to load gziped json file using xopen.

ccanalyser.utils.**get\_timing** (*task\_name*=None) → Callable

Decorator: Gets the time taken by the wrapped function

**class** ccanalyser.utils.NaturalOrderGroup (*name*=None, *commands*=None, \*\**attrs*)

Bases: click.core.Group

Simple class to ensure subcommand order is maintained by click.

**list\_commands** (*ctx*)

Returns a list of subcommand names in the order they should appear.

ccanalyser.utils.**zap\_files** (*files*)

Runs cgatcore zap\_files on all inputs

ccanalyser.utils.**get\_ucsc\_color** (*color*) → str

Converts rgb to UCSC compatable colours

ccanalyser.utils.**get\_colors** (*items*: Iterable, *colors*: Optional[Iterable] = None)

Extracts the appropriate number of colours for the items and formats them for UCSC.

ccanalyser.utils.**make\_group\_track** (*bigwigs*: list, *key*: Union[callable, str, int], *overlay*=True)

→ dict  
Generates a UCSC super track by grouping inputs by the provided key.

**class** ccanalyser.utils.PysamFakeEntry (*name*, *sequence*, *quality*)

Bases: object

Testing class used to supply a pysam FastqProxy like object

ccanalyser.utils.**convert\_to\_bedtool** (*bed*: Union[str, pybedtools.bedtool.BedTool, pandas.core.frame.DataFrame]) → pybedtools.bedtool.BedTool

Converts a str or pd.DataFrame to a pybedtools.BedTool object

ccanalyser.utils.**convert\_bed\_to\_dataframe** (*bed*: Union[str, pybedtools.bedtool.BedTool, pandas.core.frame.DataFrame]) → pandas.core.frame.DataFrame

Converts a bed like object (including paths to bed files) to a pd.DataFrame

```
ccanalyser.utils.format_coordinates(coordinates: Union[str, os.PathLike]) → pybed-
tools.bedtool.BedTool
Converts coordinates supplied in string format or a .bed file to a BedTool.
```

**Parameters** `coordinates` (`Union[str, os.PathLike]`) – Coordinates in the form chr:start-end or a path.

**Raises** `ValueError` – Inputs must be supplied in the correct format.

**Returns** BedTool object containing the required coordinates.

**Return type** BedTool

```
ccanalyser.utils.convert_interval_to_coords(interval: Union[pybedtools.cbedtools.Interval,
dict], named=False) → str
```

Converts interval object to standard genomic coordinates.

e.g. chr1:1000-2000

**Parameters** `interval` (`Union[pybedtools.Interval, dict]`) – Interval to convert.

**Returns** Genomic coordinates in the format chr:start-end

**Return type** str



---

CHAPTER  
SIX

---

## GLOSSARY

**fragment** The entire fragment of DNA that is sequenced. e.g.:

```
--R1-->
-----fragment-----
<----R2--
```

**slice** The results of *in silico* restriction digest of a fragment. e.g.:

```
-----GATC-----GATC-----
--slice1-- slice2-- slice3----
```

**capture** A slice that overlaps one of the supplied viewpoints. e.g.:

```
---capture-slice-----
-----viewpoint-----GATC-----
```

**reporter** A mapped slice from the same fragment as a capture slice used to report an interaction between two the viewpoint and another genomic region. e.g.:

```
---fragment1-capture-slice-GATC--fragment1-slice2---
```



## PYTHON MODULE INDEX

### C

ccanalyser.cli.alignments\_annotate, 25  
ccanalyser.cli.alignments\_deduplicate,  
    26  
ccanalyser.cli.alignments\_filter, 27  
ccanalyser.cli.fastq\_deduplicate, 29  
ccanalyser.cli.fastq\_digest, 30  
ccanalyser.cli.fastq\_split, 31  
ccanalyser.cli.genome\_digest, 31  
ccanalyser.cli.reporters\_count, 32  
ccanalyser.cli.reporters\_differential,  
    33  
ccanalyser.cli.reporters\_heatmap, 34  
ccanalyser.cli.reporters\_pileup, 35  
ccanalyser.cli.reporters\_store, 35  
ccanalyser.cli.tsv\_aggregate, 37  
ccanalyser.pipeline.pipeline, 37  
ccanalyser.tools.annotate, 41  
ccanalyser.tools.deduplicate, 42  
ccanalyser.tools.digest, 43  
ccanalyser.tools.filter, 45  
ccanalyser.tools.io, 53  
ccanalyser.tools.pileup, 54  
ccanalyser.tools.plotting, 56  
ccanalyser.tools.statistics, 56  
ccanalyser.tools.storage, 56  
ccanalyser.utils, 59



# INDEX

## Symbols

```
--actions <actions>
    ccanalyser-alignments-annotate
        command line option, 11
--annotations <annotations>
    ccanalyser-alignments-filter
        command line option, 14
--bam <bam>
    ccanalyser-alignments-filter
        command line option, 14
--bed_files <bed_files>
    ccanalyser-alignments-annotate
        command line option, 11
--binsize <binsize>
    ccanalyser-reporters-pileup
        command line option, 21
--binsizes <binsizes>
    ccanalyser-reporters-store-bins
        command line option, 23
--buffer <buffer>
    ccanalyser-alignments-deduplicate-identify
        command line option, 12
    ccanalyser-alignments-deduplicate-remove
        command line option, 13
--capture_name <capture_name>
    ccanalyser-reporters-differential
        command line option, 20
    ccanalyser-reporters-store-fragments
        command line option, 24
--capture_names <capture_names>
    ccanalyser-reporters-pileup
        command line option, 21
    ccanalyser-reporters-plot command
        line option, 22
--capture_viewpoints
    <capture_viewpoints>
    ccanalyser-reporters-differential
        command line option, 20
    ccanalyser-reporters-store-fragments
        command line option, 24
--cmap <cmap>
    ccanalyser-reporters-plot command
                    line option, 22
--compression_level
    <compression_level>
    ccanalyser-fastq-deduplicate-remove
        command line option, 16
    ccanalyser-fastq-digest command
        line option, 17
    ccanalyser-fastq-split command
        line option, 17
--conversion_tables
    <conversion_tables>
    ccanalyser-reporters-store-bins
        command line option, 23
--coordinates <coordinates>
    ccanalyser-reporters-plot command
        line option, 22
--custom_filtering <custom_filtering>
    ccanalyser-alignments-filter
        command line option, 14
--dense
    ccanalyser-reporters-pileup
        command line option, 21
--design_matrix <design_matrix>
    ccanalyser-reporters-differential
        command line option, 20
--duplicated_ids <duplicated_ids>
    ccanalyser-alignments-deduplicate-remove
        command line option, 13
    ccanalyser-fastq-deduplicate-remove
        command line option, 16
--duplicates <duplicates>
    ccanalyser-alignments-annotate
        command line option, 12
--fragment_map <fragment_map>
    ccanalyser-reporters-store-fragments
        command line option, 24
--genome <genome>
    ccanalyser-reporters-store-fragments
        command line option, 24
--grouping_col <grouping_col>
    ccanalyser-reporters-differential
        command line option, 20
```

```
--gzip
    ccanalyser-alignments-filter
        command line option, 14
    ccanalyser-fastq-deduplicate-remove
        command line option, 16
    ccanalyser-fastq-split command
        line option, 17
    ccanalyser-reporters-pileup
        command line option, 21
--help
    ccanalyser-pipeline command line
        option, 19
--invalid_bed_action
    <invalid_bed_action>
    ccanalyser-alignments-annotate
        command line option, 12
--keep_cutsite <keep_cutsite>
    ccanalyser-fastq-digest command
        line option, 17
--logfile <logfile>
    ccanalyser-genome-digest command
        line option, 18
--method <method>
    ccanalyser-fastq-split command
        line option, 17
--minimum_slice_length
    <minimum_slice_length>
    ccanalyser-fastq-digest command
        line option, 17
--mode <mode>
    ccanalyser-fastq-digest command
        line option, 17
--n_cores <n_cores>
    ccanalyser-alignments-annotate
        command line option, 12
    ccanalyser-fastq-digest command
        line option, 17
    ccanalyser-reporters-store-bins
        command line option, 23
--n_reads <n_reads>
    ccanalyser-fastq-split command
        line option, 17
--names <names>
    ccanalyser-alignments-annotate
        command line option, 11
--no-gzip
    ccanalyser-alignments-filter
        command line option, 14
    ccanalyser-fastq-deduplicate-remove
        command line option, 16
    ccanalyser-fastq-split command
        line option, 17
--normalisation <normalisation>
    ccanalyser-reporters-plot command
        line option, 22
--overlap_fraction <overlap_fraction>
    ccanalyser-reporters-store-bins
        command line option, 23
--overlap_fractions
    <overlap_fractions>
    ccanalyser-alignments-annotate
        command line option, 11
--read_buffer <read_buffer>
    ccanalyser-fastq-deduplicate-parse
        command line option, 15
    ccanalyser-fastq-deduplicate-remove
        command line option, 16
```

```

        command line option, 16
    ccanalyser-fastq-digest command
        line option, 17
--read_type <read_type>
    ccanalyser-alignments-deduplicate-identify
        command line option, 12
    ccanalyser-alignments-deduplicate-remove
        command line option, 13
    ccanalyser-alignments-filter
        command line option, 14
--recognition_site <recognition_site>
    ccanalyser-genome-digest command
        line option, 18
--remove_capture
    ccanalyser-reporters-count command
        line option, 19
    ccanalyser-reporters-plot command
        line option, 22
--remove_cutsite <remove_cutsite>
    ccanalyser-genome-digest command
        line option, 18
--remove_exclusions
    ccanalyser-reporters-count command
        line option, 19
--resolution <resolution>
    ccanalyser-reporters-plot command
        line option, 22
--restriction_enzyme
    <restriction_enzyme>
    ccanalyser-fastq-digest command
        line option, 17
--sample_name <sample_name>
    ccanalyser-alignments-deduplicate-remove
        command line option, 13
    ccanalyser-alignments-filter
        command line option, 14
    ccanalyser-fastq-deduplicate-remove
        command line option, 16
    ccanalyser-fastq-digest command
        line option, 17
--scale_factor <scale_factor>
    ccanalyser-reporters-pileup
        command line option, 21
    ccanalyser-reporters-store-bins
        command line option, 23
--sort
    ccanalyser-genome-digest command
        line option, 18
--sparse
    ccanalyser-reporters-pileup
        command line option, 21
--stats_prefix <stats_prefix>
    ccanalyser-alignments-deduplicate-remove
    ccanalyser-alignments-filter
        command line option, 13
ccanalyser-alignments-filter
    command line option, 14
ccanalyser-fastq-deduplicate-remove
    command line option, 16
ccanalyser-fastq-digest command
    line option, 17
<subsample>
    ccanalyser-reporters-count command
        line option, 19
--suffix <suffix>
    ccanalyser-reporters-store-fragments
        command line option, 24
--threshold_count <threshold_count>
    ccanalyser-reporters-differential
        command line option, 20
--threshold_mean <threshold_mean>
    ccanalyser-reporters-differential
        command line option, 20
--threshold_q <threshold_q>
    ccanalyser-reporters-differential
        command line option, 20
--vmax <vmax>
    ccanalyser-reporters-plot command
        line option, 22
--vmin <vmin>
    ccanalyser-reporters-plot command
        line option, 22
-a
    ccanalyser-alignments-annotate
        command line option, 11
    ccanalyser-alignments-filter
        command line option, 14
ccanalyser-alignments-annotate
    command line option, 11
ccanalyser-alignments-filter
    command line option, 14
ccanalyser-reporters-store-bins
    command line option, 23
-c
    ccanalyser-reporters-differential
        command line option, 20
    ccanalyser-reporters-plot command
        line option, 22
    ccanalyser-reporters-store-fragments
        command line option, 24
-d
    ccanalyser-alignments-deduplicate-remove
        command line option, 13
    ccanalyser-fastq-deduplicate-remove
        command line option, 16
-f
    ccanalyser-alignments-annotate
        command line option, 11

```

ccanalyser-reporters-store-fragments  
    command line option, 24

-g  
    ccanalyser-reporters-store-fragments  
        command line option, 24

-h  
    ccanalyser-pipeline command line  
        option, 19

-l  
    ccanalyser-genome-digest command  
        line option, 18

-m  
    ccanalyser-fastq-digest command  
        line option, 17

    ccanalyser-fastq-split command  
        line option, 17

-n  
    ccanalyser-alignments-annotate  
        command line option, 11

    ccanalyser-fastq-split command  
        line option, 17

    ccanalyser-reporters-differential  
        command line option, 20

    ccanalyser-reporters-pileup  
        command line option, 21

    ccanalyser-reporters-plot command  
        line option, 22

    ccanalyser-reporters-store-fragments  
        command line option, 24

-o  
    ccanalyser-alignments-annotate  
        command line option, 12

    ccanalyser-alignments-deduplicate-identify  
        command line option, 12

    ccanalyser-alignments-deduplicate-remove  
        command line option, 13

    ccanalyser-alignments-filter  
        command line option, 14

    ccanalyser-fastq-deduplicate-identify  
        command line option, 15

    ccanalyser-fastq-deduplicate-parse  
        command line option, 15

    ccanalyser-fastq-deduplicate-remove  
        command line option, 16

    ccanalyser-fastq-digest command  
        line option, 17

    ccanalyser-fastq-split command  
        line option, 17

    ccanalyser-genome-digest command  
        line option, 18

    ccanalyser-reporters-count command  
        line option, 19

    ccanalyser-reporters-differential  
        command line option, 20

ccanalyser-reporters-pileup  
    command line option, 21

ccanalyser-reporters-plot command  
    line option, 22

ccanalyser-reporters-store-bins  
    command line option, 23

ccanalyser-reporters-store-fragments  
    command line option, 24

ccanalyser-reporters-store-merge  
    command line option, 24

-p  
    ccanalyser-alignments-annotate  
        command line option, 12

    ccanalyser-fastq-digest command  
        line option, 17

    ccanalyser-reporters-store-bins  
        command line option, 23

-r  
    ccanalyser-fastq-digest command  
        line option, 17

    ccanalyser-genome-digest command  
        line option, 18

    ccanalyser-reporters-plot command  
        line option, 22

**A**

alignments\_deduplicate\_collate() (in module `ccanalyser.pipeline.pipeline`), 39

alignments\_deduplicate\_fragments() (in module `ccanalyser.pipeline.pipeline`), 39

alignments\_deduplicate\_slices() (in module `ccanalyser.pipeline.pipeline`), 39

alignments\_filter() (in module `ccanalyser.pipeline.pipeline`), 39

alignments\_index() (in module `ccanalyser.pipeline.pipeline`), 39

alignments\_merge() (in module `ccanalyser.pipeline.pipeline`), 39

annotate() (in module `ccanalyser.cli.alignments_annotate`), 25

annotate\_alignments() (in module `ccanalyser.pipeline.pipeline`), 39

annotate\_make\_exclusion\_bed() (in module `ccanalyser.pipeline.pipeline`), 39

annotate\_sort\_blacklist() (in module `ccanalyser.pipeline.pipeline`), 39

annotate\_sort\_viewpoints() (in module `ccanalyser.pipeline.pipeline`), 39

**B**

bed1 (`ccanalyser.tools.annotate.BedIntersection` attribute), 41

bed2 (`ccanalyser.tools.annotate.BedIntersection` attribute), 41

bed\_has\_duplicate\_names () (*in module ccanalyser.utils*), 59  
 bed\_has\_name () (*in module ccanalyser.utils*), 59  
 bedgraph () (*ccanalyser.tools.pileup.CoolerBedGraph property*), 54  
 bedgraph () (*in module ccanalyser.cli.reporters\_pileup*), 35  
 BedIntersection (*class in ccanalyser.tools.annotate*), 41  
 bin\_conversion\_table () (*ccanalyser.tools.storage.CoolerBinner property*), 58  
 bin\_conversion\_table () (*ccanalyser.tools.storage.GenomicBinner property*), 58  
 binner (*ccanalyser.tools.storage.CoolerBinner attribute*), 58  
 bins () (*ccanalyser.tools.storage.CoolerBinner property*), 58  
 bins () (*ccanalyser.tools.storage.GenomicBinner property*), 57  
 bins () (*in module ccanalyser.cli.reporters\_store*), 36  
 binsize (*ccanalyser.tools.storage.CoolerBinner attribute*), 58  
 binsize (*ccanalyser.tools.storage.GenomicBinner attribute*), 57

## C

capture, 63  
 capture\_bins () (*ccanalyser.tools.storage.CoolerBinner property*), 58  
 capture\_name (*ccanalyser.tools.pileup.CoolerBedGraph attribute*), 54  
 capture\_site\_stats () (*ccanalyser.tools.filter.CCSliceFilter property*), 49  
 captures () (*ccanalyser.tools.filter.CCSliceFilter property*), 49  
 ccanalyser.cli.alignments\_annotation module, 25  
 ccanalyser.cli.alignments\_deduplicate module, 26  
 ccanalyser.cli.alignments\_filter module, 27  
 ccanalyser.cli.fastq\_deduplicate module, 29  
 ccanalyser.cli.fastq\_digest module, 30  
 ccanalyser.cli.fastq\_split module, 31  
 ccanalyser.cli.genome\_digest module, 31  
 ccanalyser.cli.reporters\_count module, 32  
 ccanalyser.cli.reporters\_differential module, 33  
 ccanalyser.cli.reporters\_heatmap module, 34  
 ccanalyser.cli.reporters\_pileup module, 35  
 ccanalyser.cli.reporters\_store module, 35  
 ccanalyser.cli.tsv\_aggregate module, 37  
 ccanalyser.pipeline.pipeline module, 37  
 ccanalyser.tools.annotate module, 41  
 ccanalyser.tools.deduplicate module, 42  
 ccanalyser.tools.digest module, 43  
 ccanalyser.tools.filter module, 45  
 ccanalyser.tools.io module, 53  
 ccanalyser.tools.pileup module, 54  
 ccanalyser.tools.plotting module, 56  
 ccanalyser.tools.statistics module, 56  
 ccanalyser.tools.storage module, 56  
 ccanalyser.utils module, 59

ccanalyser-alignments-annotate command line option  
 --actions <actions>, 11  
 --bed\_files <bed\_files>, 11  
 --duplicates <duplicates>, 12  
 --invalid\_bed\_action <invalid\_bed\_action>, 12  
 --n\_cores <n\_cores>, 12  
 --names <names>, 11  
 --output <output>, 12  
 --overlap\_fractions <overlap\_fractions>, 11

-a, 11  
 -b, 11  
 -f, 11  
 -n, 11  
 -o, 12  
 -p, 12  
 SLICES, 12

ccanalyser-alignments-deduplicate-identify command line option

```
--buffer <buffer>, 12
--output <output>, 12
--read_type <read_type>, 12
-o, 12
FRAGMENTS_FN, 13
ccanalyser-alignments-deduplicate-removeccanalyser-fastq-digest command line
    command line option
--buffer <buffer>, 13
--duplicated_ids <duplicated_ids>, 13
--output <output>, 13
--read_type <read_type>, 13
--sample_name <sample_name>, 13
--stats_prefix <stats_prefix>, 13
-d, 13
-o, 13
SLICES_FN, 13
ccanalyser-alignments-filter command
    line option
--annotations <annotations>, 14
--bam <bam>, 14
--custom_filtering
    <custom_filtering>, 14
--gzip, 14
--no-gzip, 14
--output_prefix <output_prefix>, 14
--read_type <read_type>, 14
--sample_name <sample_name>, 14
--stats_prefix <stats_prefix>, 14
-a, 14
-b, 14
-o, 14
METHOD, 14
ccanalyser-fastq-deduplicate-identify
    command line option
--output <output>, 15
-o, 15
INPUT_FILES, 15
ccanalyser-fastq-deduplicate-parse
    command line option
--output <output>, 15
--read_buffer <read_buffer>, 15
-o, 15
INPUT_FILES, 15
ccanalyser-fastq-deduplicate-remove
    command line option
--compression_level
    <compression_level>, 16
--duplicated_ids <duplicated_ids>, 16
--gzip, 16
--no-gzip, 16
--output_prefix <output_prefix>, 16
--read_buffer <read_buffer>, 16
--sample_name <sample_name>, 16
--stats_prefix <stats_prefix>, 16
-d, 16
-o, 16
INPUT_FILES, 16
ccanalyser-pipeline command line
    option
--help, 19
```

```

-h, 19
MODE, 19
PIPELINE_OPTIONS, 19
ccanalyser-reporters-count command
    line option
    --output <output>, 19
    --remove_capture, 19
    --remove_exclusions, 19
    --subsample <subsample>, 19
    -o, 19
    REPORTERS, 20
ccanalyser-reporters-differential
    command line option
    --capture_name <capture_name>, 20
    --capture_viewpoints
        <capture_viewpoints>, 20
    --design_matrix <design_matrix>, 20
    --grouping_col <grouping_col>, 20
    --output_prefix <output_prefix>, 20
    --threshold_count
        <threshold_count>, 20
    --threshold_mean <threshold_mean>,
        20
    --threshold_q <threshold_q>, 20
    -c, 20
    -n, 20
    -o, 20
    UNION_BEDGRAPH, 21
ccanalyser-reporters-pileup command
    line option
    --binsize <binsize>, 21
    --capture_names <capture_names>, 21
    --dense, 21
    --gzip, 21
    --normalise, 21
    --output_prefix <output_prefix>, 21
    --scale_factor <scale_factor>, 21
    --sparse, 21
    -n, 21
    -o, 21
    COOLER_FN, 21
ccanalyser-reporters-plot command line
    option
    --capture_names <capture_names>, 22
    --cmap <cmap>, 22
    --coordinates <coordinates>, 22
    --normalisation <normalisation>, 22
    --output_prefix <output_prefix>, 22
    --remove_capture, 22
    --resolution <resolution>, 22
    --vmax <vmax>, 22
    --vmin <vmin>, 22
    -c, 22
    -n, 22
-o, 22
-r, 22
COOLER_FN, 22
ccanalyser-reporters-store-bins
    command line option
    --binsizes <binsizes>, 23
    --conversion_tables
        <conversion_tables>, 23
    --n_cores <n_cores>, 23
    --normalise, 23
    --output <output>, 23
    --overlap_fraction
        <overlap_fraction>, 23
    --scale_factor <scale_factor>, 23
-b, 23
-o, 23
-p, 23
COOLER_FN, 23
ccanalyser-reporters-store-fragments
    command line option
    --capture_name <capture_name>, 24
    --capture_viewpoints
        <capture_viewpoints>, 24
    --fragment_map <fragment_map>, 24
    --genome <genome>, 24
    --output <output>, 24
    --suffix <suffix>, 24
    -c, 24
    -f, 24
    -g, 24
    -n, 24
    -o, 24
    COUNTS, 24
ccanalyser-reporters-store-merge
    command line option
    --output <output>, 24
    -o, 24
    COOLERS, 24
CCBedgraph (class in ccanalyser.tools.pileup), 55
CCMatrix (class in ccanalyser.tools.plotting), 56
CCSliceFilter (class in ccanalyser.tools.filter), 47
chrom (ccanalyser.tools.digest.DigestedChrom attribute), 43
chromsizes (ccanalyser.tools.storage.GenomicBinner attribute), 57
cis_or_trans_stats() (ccanalyser.tools.filter.CCSliceFilter property), 49
cis_or_trans_stats() (ccanalyser.tools.filter.TiledCSliceFilter property), 52
collate_cis_trans_data() (in module ccanalyser.tools.statistics), 56
collate_histogram_data() (in module ccanalyser.tools.statistics), 56

```

collate\_read\_data() (in module `ccanalyser.tools.statistics`), 56  
collate\_slice\_data() (in module `ccanalyser.tools.statistics`), 56  
collate\_statistics() (in module `ccanalyser.cli.fastq_digest`), 30  
concat\_tsvs() (in module `ccanalyser.cli.tsv_aggregate`), 37  
convert\_bed\_to\_dataframe() (in module `ccanalyser.utils`), 60  
convert\_interval\_to\_coords() (in module `ccanalyser.utils`), 61  
convert\_to\_bedtool() (in module `ccanalyser.utils`), 60  
cooler (`ccanalyser.tools.pileup.CoolerBedGraph` attribute), 54  
cooler (`ccanalyser.tools.storage.CoolerBinner` attribute), 58  
COOLER\_FN  
    ccanalyser-reporters-pileup command line option, 21  
    ccanalyser-reporters-plot command line option, 22  
    ccanalyser-reporters-store-bins command line option, 23  
CoolerBedGraph (class in `ccanalyser.tools.pileup`), 54  
CoolerBedGraphWindowed (class in `ccanalyser.tools.pileup`), 55  
CoolerBinner (class in `ccanalyser.tools.storage`), 58  
COOLERS  
    ccanalyser-reporters-store-merge command line option, 24  
coordinates() (ccanalyser.tools.pileup.CCBedgraph property), 55  
count() (in module `ccanalyser.cli.reporters_count`), 32  
count\_re\_site\_combinations() (in module `ccanalyser.cli.reporters_count`), 32  
COUNTS  
    ccanalyser-reporters-store-fragments command line option, 24  
create\_cooler\_cc() (in module `ccanalyser.tools.storage`), 56  
cycle\_argument() (in module `ccanalyser.cli.alignments_annotate`), 25

**D**

DeduplicationStatistics (class in `ccanalyser.tools.statistics`), 56  
df() (`ccanalyser.tools.statistics.DeduplicationStatistics` property), 56  
differential() (in module `ccanalyser.cli.reporters_differential`), 33

digest() (in module `ccanalyser.cli.fastq_digest`), 30  
digest() (in module `ccanalyser.cli.genome_digest`), 31  
DigestedChrom (class in `ccanalyser.tools.digest`), 43  
DigestedRead (class in `ccanalyser.tools.digest`), 44  
DigestionStats (class in `ccanalyser.tools.statistics`), 56  
duplicated\_ids (ccanalyser.tools.deduplicate.ReadDuplicateRemovalProcess attribute), 42

**E**

extract\_trimming\_stats() (in module `ccanalyser.tools.statistics`), 56

**F**

fastq\_alignment() (in module `ccanalyser.pipeline.pipeline`), 39  
fastq\_digest\_combined() (in module `ccanalyser.pipeline.pipeline`), 38  
fastq\_digest\_non\_combined() (in module `ccanalyser.pipeline.pipeline`), 39  
fastq\_duplicates\_identify() (in module `ccanalyser.pipeline.pipeline`), 38  
fastq\_duplicates\_parse() (in module `ccanalyser.pipeline.pipeline`), 38  
fastq\_duplicates\_remove() (in module `ccanalyser.pipeline.pipeline`), 38  
fastq\_flash() (in module `ccanalyser.pipeline.pipeline`), 38  
fastq\_multiqc() (in module `ccanalyser.pipeline.pipeline`), 38  
fastq\_preprocessing() (in module `ccanalyser.pipeline.pipeline`), 39  
fastq\_qc() (in module `ccanalyser.pipeline.pipeline`), 38  
fastq\_split() (in module `ccanalyser.pipeline.pipeline`), 38  
fastq\_trim() (in module `ccanalyser.pipeline.pipeline`), 38  
FastqReaderProcess (class in `ccanalyser.tools.io`), 53  
FastqReadFormatterProcess (class in `ccanalyser.tools.io`), 53  
FastqWriterProcess (class in `ccanalyser.tools.io`), 53  
FastqWriterSplitterProcess (class in `ccanalyser.tools.io`), 53  
filter() (in module `ccanalyser.cli.alignments_filter`), 27  
filter\_slices() (`ccanalyser.tools.filter.SliceFilter` method), 46  
filter\_stages (`ccanalyser.tools.filter.CCSliceFilter` attribute), 48

filter\_stages (*ccanalyser.tools.filter.SliceFilter attribute*), 45  
 filter\_stages (*ccanalyser.tools.filter.TiledCSliceFilter attribute*), 51  
 filter\_stages (*ccanalyser.tools.filter.TriCSliceFilter attribute*), 50  
 filter\_stats (*ccanalyser.tools.filter.CCSliceFilter attribute*), 48  
 filter\_stats (*ccanalyser.tools.filter.SliceFilter attribute*), 46  
 filter\_stats (*ccanalyser.tools.filter.TiledCSliceFilter attribute*), 52  
 filter\_stats (*ccanalyser.tools.filter.TriCSliceFilter attribute*), 51  
 filter\_stats () (*ccanalyser.tools.filter.SliceFilter property*), 46  
 filtered (*ccanalyser.tools.statistics.DigestionStats attribute*), 56  
 filters () (*ccanalyser.tools.filter.SliceFilter property*), 46  
 format\_coordinates () (in module *ccanalyser.utils*), 60  
 format\_index\_var () (in module *ccanalyser.cli.tsv\_aggregate*), 37  
 frag\_stats () (*ccanalyser.tools.filter.CCSliceFilter property*), 49  
 fragment, 63  
 fragment\_indexes (*ccanalyser.tools.digest.DigestedChrom attribute*), 43  
 fragment\_min\_len (*ccanalyser.tools.digest.DigestedChrom attribute*), 43  
 fragment\_number\_offset (*ccanalyser.tools.digest.DigestedChrom attribute*), 43  
 fragments (*ccanalyser.tools.filter.CCSliceFilter attribute*), 48  
 fragments (*ccanalyser.tools.filter.SliceFilter attribute*), 45  
 fragments (*ccanalyser.tools.filter.TiledCSliceFilter attribute*), 51  
 fragments (*ccanalyser.tools.filter.TriCSliceFilter attribute*), 50  
 fragments (*ccanalyser.tools.storage.GenomicBinner attribute*), 57  
 fragments () (*ccanalyser.tools.digest.DigestedChrom property*), 44  
 fragments () (*ccanalyser.tools.filter.CCSliceFilter property*), 48  
 fragments () (*ccanalyser.tools.filter.SliceFilter prop-  
 erty*), 46  
 fragments () (*ccanalyser.tools.filter.TiledCSliceFilter property*), 52  
 fragments () (in module *ccanalyser.cli.reporters\_store*), 35  
 FRAGMENTS\_FN  
 ccanalyser-alignments-deduplicate-identify command line option, 13  
 full () (in module *ccanalyser.pipeline.pipeline*), 41

## G

generate\_bin\_conversion\_tables () (in module *ccanalyser.pipeline.pipeline*), 40  
 genome\_digest () (in module *ccanalyser.pipeline.pipeline*), 38  
 GenomicBinner (class in *ccanalyser.tools.storage*), 57  
 get\_capture\_bins () (in module *ccanalyser.tools.storage*), 56  
 get\_capture\_coords () (in module *ccanalyser.tools.storage*), 56  
 get\_chromosome\_from\_name () (in module *ccanalyser.cli.reporters\_differential*), 33  
 get\_colors () (in module *ccanalyser.utils*), 60  
 get\_human\_readable\_number\_of\_bp () (in module *ccanalyser.utils*), 59  
 get\_matrix () (*ccanalyser.tools.plotting.CCMATRIX method*), 56  
 get\_matrix\_normalised () (*ccanalyser.tools.plotting.CCMATRIX method*), 56  
 get\_re\_site () (in module *ccanalyser.utils*), 59  
 get\_recognition\_site\_indexes () (*ccanalyser.tools.digest.DigestedChrom method*), 43  
 get\_recognition\_site\_indexes () (*ccanalyser.tools.digest.DigestedRead method*), 45  
 get\_timing () (in module *ccanalyser.utils*), 60  
 get\_ucsc\_color () (in module *ccanalyser.utils*), 60  
 get\_unfiltered\_slices () (*ccanalyser.tools.filter.SliceFilter method*), 46

## H

has\_slices (*ccanalyser.tools.digest.DigestedRead attribute*), 44  
 hash\_column () (in module *ccanalyser.utils*), 60  
 hash\_seed (*ccanalyser.tools.deduplicate.ReadDeduplicationParserProcess attribute*), 42  
 hash\_seed (*ccanalyser.tools.deduplicate.ReadDuplicateRemovalProcess attribute*), 43  
 hub\_make () (in module *ccanalyser.pipeline.pipeline*), 40  
 hub\_write\_path () (in module *ccanalyser.pipeline.pipeline*), 40

|

identify() (in module `ccanalyser.cli.alignments_deduplicate`), 26  
identify() (in module `ccanalyser.cli.fastaq_deduplicate`), 29  
identify\_differential\_interactions() (in module `ccanalyser.pipeline.pipeline`), 40  
INPUT\_FASTA  
    ccanalyser-genome-digest command line option, 18  
INPUT\_FASTQ  
    ccanalyser-fastq-digest command line option, 17  
input\_file (`ccanalyser.tools.io.FastqReaderProcess` attribute), 53  
INPUT\_FILES  
    ccanalyser-fastq-deduplicate-identify METHOD  
        command line option, 15  
    ccanalyser-fastq-deduplicate-parse  
        command line option, 15  
    ccanalyser-fastq-deduplicate-remove  
        command line option, 16  
    ccanalyser-fastq-split command line option, 18  
inq(`ccanalyser.tools.deduplicate.ReadDeduplicationParserProcess` attribute), 42  
inq(`ccanalyser.tools.deduplicate.ReadDuplicateRemovalProcess` attribute), 42  
intersect\_bins() (in module `ccanalyser.utils`), 60  
intersection() (ccanalyser.tools.annotate.BedIntersection property), 41  
intersection\_name (ccanalyser.tools.annotate.BedIntersection attribute), 41  
invalid\_bed\_action (ccanalyser.tools.annotate.BedIntersection attribute), 41  
invert\_dict() (in module `ccanalyser.utils`), 59  
is\_compressed() (in module `ccanalyser.cli.tsv_aggregate`), 37  
is\_none() (in module `ccanalyser.utils`), 59  
is\_off() (in module `ccanalyser.utils`), 59  
is\_on() (in module `ccanalyser.utils`), 59  
is\_valid\_bed() (in module `ccanalyser.utils`), 59

J

join\_tsvs() (in module `ccanalyser.cli.tsv_aggregate`), 37

L

link\_bins() (in module `ccanalyser.tools.storage`), 59  
list\_commands() (ccanalyser.utils.NaturalOrderGroup method), 60

load\_json() (in module `ccanalyser.utils`), 60  
load\_tsv() (in module `ccanalyser.cli.tsv_aggregate`), 37

M

main() (in module `ccanalyser.cli.tsv_aggregate`), 37  
make\_group\_track() (in module `ccanalyser.utils`), 60  
merge() (in module `ccanalyser.cli.reporters_store`), 36  
merge\_annotations() (in module `ccanalyser.cli.alignments_filter`), 27  
merge\_dictionaries() (in module `ccanalyser.utils`), 59  
merged\_captures\_and\_reporters() (ccanalyser.tools.filter.CCSliceFilter property), 49

METHOD

ccanalyser-alignments-filter command line option, 14  
min\_frac (`ccanalyser.tools.annotate.BedIntersection` attribute), 41  
min\_overlap (ccanalyser.tools.storage.GenomicBinner attribute), 57  
slice\_len (ccanalyser.tools.digest.DigestedRead attribute), 44

MODE

ccanalyser-pipeline command line option, 19  
modify\_pipeline\_params\_dict() (in module `ccanalyser.pipeline.pipeline`), 38

module

ccanalyser.cli.alignments\_annotation, 25  
ccanalyser.cli.alignments\_deduplicate, 26  
ccanalyser.cli.alignments\_filter, 27  
ccanalyser.cli.fastq\_deduplicate, 29  
ccanalyser.cli.fastq\_digest, 30  
ccanalyser.cli.fastq\_split, 31  
ccanalyser.cli.genome\_digest, 31  
ccanalyser.cli.reporters\_count, 32  
ccanalyser.cli.reporters\_differential, 33  
ccanalyser.cli.reporters\_heatmap, 34  
ccanalyser.cli.reporters\_pileup, 35  
ccanalyser.cli.reporters\_store, 35  
ccanalyser.cli.tsv\_aggregate, 37  
ccanalyser.pipeline.pipeline, 37  
ccanalyser.tools.annotate, 41  
ccanalyser.tools.deduplicate, 42  
ccanalyser.tools.digest, 43  
ccanalyser.tools.filter, 45

ccanalyser.tools.io, 53  
 ccanalyser.tools.pileup, 54  
 ccanalyser.tools.plotting, 56  
 ccanalyser.tools.statistics, 56  
 ccanalyser.tools.storage, 56  
 ccanalyser.utils, 59

**N**

n\_cis\_interactions (ccanalyser.tools.storage.CoolerBinner attribute), 58  
 n\_cores (ccanalyser.tools.annotate.BedIntersection attribute), 41  
 n\_cores (ccanalyser.tools.storage.CoolerBinner attribute), 58  
 n\_cores (ccanalyser.tools.storage.GenomicBinner attribute), 57  
 n\_subprocesses (ccanalyser.tools.io.FastqReaderProcess attribute), 53  
 NaturalOrderGroup (class in ccanalyser.utils), 60  
 normalise\_bedgraph () (ccanalyser.tools.pileup.CoolerBedGraph method), 55  
 normalise\_bedgraph () (ccanalyser.tools.pileup.CoolerBedGraphWindowed method), 55  
 normalise\_pixels () (ccanalyser.tools.storage.CoolerBinner method), 58

**O**

only\_cis (ccanalyser.tools.pileup.CoolerBedGraph attribute), 54  
 open\_logfile () (in module ccanalyser.utils), 59  
 outq (ccanalyser.tools.deduplicate.ReadDeduplicationParserProcess attribute), 42  
 outq (ccanalyser.tools.deduplicate.ReadDuplicateRemovalProcess attribute), 42  
 outq (ccanalyser.tools.io.FastqReaderProcess attribute), 53

**P**

parse () (in module ccanalyser.cli.fastq\_deduplicate), 29  
 parse\_alignment () (in module ccanalyser.tools.io), 54  
 parse\_bam () (in module ccanalyser.tools.io), 54  
 parse\_chromosomes () (in module ccanalyser.cli.genome\_digest), 31  
 pipeline\_make\_report () (in module ccanalyser.pipeline.pipeline), 40  
 pipeline\_merge\_stats () (in module ccanalyser.pipeline.pipeline), 40

PIPELINE\_OPTIONS  
 ccanalyser-pipeline command line option, 19  
 pixel\_conversion\_table () (ccanalyser.tools.storage.CoolerBinner property), 58  
 pixels () (ccanalyser.tools.storage.CoolerBinner property), 58  
 plot () (in module ccanalyser.cli.reporters\_heatmap), 34  
 plot\_matrix () (in module ccanalyser.cli.reporters\_heatmap), 34  
 post\_annotation () (in module ccanalyser.pipeline.pipeline), 39  
 post\_ccanalyser\_analysis () (in module ccanalyser.pipeline.pipeline), 39  
 pre\_annotation () (in module ccanalyser.pipeline.pipeline), 39  
 PysamFakeEntry (class in ccanalyser.utils), 60

**R**

read (ccanalyser.tools.digest.DigestedRead attribute), 44  
 read\_buffer (ccanalyser.tools.io.FastqReaderProcess attribute), 53  
 read\_counter (ccanalyser.tools.io.FastqReaderProcess attribute), 53  
 read\_number (ccanalyser.tools.statistics.DigestionStats attribute), 56  
 read\_stats (ccanalyser.tools.filter.CCSliceFilter attribute), 48  
 read\_stats (ccanalyser.tools.filter.SliceFilter attribute), 45  
 read\_stats (ccanalyser.tools.filter.TiledCSliceFilter attribute), 52  
 read\_stats (ccanalyser.tools.filter.TriCSliceFilter attribute), 51  
 read\_stats () (ccanalyser.tools.filter.SliceFilter property), 46  
 read\_type (ccanalyser.tools.statistics.DigestionStats attribute), 56  
 ReadDeduplicationParserProcess (class in ccanalyser.tools.deduplicate), 42  
 ReadDigestionProcess (class in ccanalyser.tools.digest), 45  
 ReadDuplicateRemovalProcess (class in ccanalyser.tools.deduplicate), 42  
 reads\_total (ccanalyser.tools.deduplicate.ReadDuplicateRemovalProcess attribute), 43

reads\_unique (ccanal-  
yser.tools.deduplicate.ReadDuplicateRemovalProcess.replace\_na() (in module ccanal-  
attribute), 43  
recognition\_len (ccanal-  
yser.tools.digest.DigestedChrom  
43  
recognition\_len (ccanal-  
yser.tools.digest.DigestedRead  
44  
recognition\_seq (ccanal-  
yser.tools.digest.DigestedChrom  
43  
recognition\_seq (ccanal-  
yser.tools.digest.DigestedRead  
44  
remove() (in module ccanal-  
yser.cli.alignments\_deduplicate), 26  
remove() (in module ccanalyser.cli.fastq\_deduplicate),  
29  
remove\_blacklisted\_slices() (ccanal-  
yser.tools.filter.SliceFilter method), 47  
remove\_dual\_capture\_fragments() (ccanal-  
yser.tools.filter.TiledCSliceFilter  
52  
remove\_duplicate\_re frags() (ccanal-  
yser.tools.filter.SliceFilter method), 47  
remove\_duplicate\_slices() (ccanal-  
yser.tools.filter.SliceFilter method), 47  
remove\_duplicate\_slices\_pe() (ccanal-  
yser.tools.filter.SliceFilter method), 47  
remove\_duplicates\_from\_bed() (in module ccan-  
alyser.cli.alignments\_annotate), 25  
remove\_excluded\_slices() (ccanal-  
yser.tools.filter.SliceFilter method), 47  
remove\_multi\_capture\_fragments() (ccanal-  
yser.tools.filter.CCSliceFilter method), 49  
remove\_multicapture\_reporters() (ccanal-  
yser.tools.filter.CCSliceFilter method), 49  
remove\_non\_capture\_fragments() (ccanal-  
yser.tools.filter.TiledCSliceFilter  
52  
remove\_non\_reporter\_fragments() (ccanal-  
yser.tools.filter.CCSliceFilter method), 49  
remove\_orphan\_slices() (ccanal-  
yser.tools.filter.SliceFilter method), 47  
remove\_slices\_outside\_capture() (ccanal-  
yser.tools.filter.TiledCSliceFilter  
52  
remove\_slices\_with\_one\_reporter() (cc-  
analyser.tools.filter.TriCSliceFilter  
method), 51  
remove\_slices\_without\_re\_frag\_assigned() run()  
(ccanalyser.tools.filter.SliceFilter method), 47  
remove\_unmapped\_slices() (ccanal-  
yser.tools.filter.SliceFilter method), 47  
replace\_na() (in module ccanal-  
yser.cli.tsv\_aggregate), 37  
reporter, 63  
REPORTERS  
ccanalyser-reporters-count command  
line option, 20  
reporters (ccanalyser.tools.filter.CCSliceFilter  
attribute), 48  
reporters (ccanalyser.tools.filter.SliceFilter  
attribute), 45  
reporters (ccanalyser.tools.filter.TiledCSliceFilter  
attribute), 51  
reporters (ccanalyser.tools.filter.TriCSliceFilter  
attribute), 50  
reporters() (ccanalyser.tools.filter.CCSliceFilter  
property), 49  
reporters() (ccanalyser.tools.filter.SliceFilter prop-  
erty), 46  
reporters() (ccanal-  
yser.tools.pileup.CoolerBedGraph  
property),  
55  
reporters\_binned() (ccanal-  
yser.tools.pileup.CoolerBedGraphWindowed  
property), 55  
reporters\_collate() (in module ccanal-  
yser.pipeline.pipeline), 39  
reporters\_count() (in module ccanal-  
yser.pipeline.pipeline), 39  
reporters\_make\_bedgraph() (in module ccanal-  
yser.pipeline.pipeline), 40  
reporters\_make\_bedgraph\_normalised() (in  
module ccanalyser.pipeline.pipeline), 40  
reporters\_make\_bigwig() (in module ccanal-  
yser.pipeline.pipeline), 40  
reporters\_make\_subtraction\_bedgraph()  
(in module ccanalyser.pipeline.pipeline), 40  
reporters\_make\_union\_bedgraph() (in mod-  
ule ccanalyser.pipeline.pipeline), 40  
reporters\_plot\_heatmap() (in module ccanal-  
yser.pipeline.pipeline), 40  
reporters\_store\_binned() (in module ccanal-  
yser.pipeline.pipeline), 40  
reporters\_store\_merged() (in module ccanal-  
yser.pipeline.pipeline), 40  
reporters\_store\_restriction\_fragment()  
(in module ccanalyser.pipeline.pipeline), 40  
run() (ccanalyser.tools.deduplicate.ReadDeduplicationParserProcess  
method), 42  
run() (ccanalyser.tools.deduplicate.ReadDuplicateRemovalProcess  
method), 43  
run() (ccanalyser.tools.digest.ReadDigestionProcess  
method), 45  
run() (ccanalyser.tools.io.FastqReaderProcess

```

        method), 53
run() (ccanalyser.tools.io.FastqReadFormatterProcess
       method), 53
run() (ccanalyser.tools.io.FastqWriterProcess method),
       54
run() (ccanalyser.tools.io.FastqWriterSplitterProcess
       method), 53
run_unix_split() (in module ccanal-
                  yser.cli.fastq_split), 31

S
save_hash_dict_path (ccanal-
                     yser.tools.deduplicate.ReadDeduplicationParserProcess
                     attribute), 42
scale_factor (ccanal-
              yser.tools.storage.CoolerBinner
              attribute), 58
score() (ccanalyser.tools.pileup.CCBedgraph
         property), 55
set_up_chromsizes() (in module ccanal-
                      yser.pipeline.pipeline), 38
slice, 63
slice_indexes (ccanal-
                 yser.tools.digest.DigestedRead
                 attribute), 44
slice_number_offset (ccanal-
                     yser.tools.digest.DigestedRead
                     attribute), 44
slice_stats (ccanalyser.tools.filter.CCSliceFilter
            attribute), 48
slice_stats (ccanalyser.tools.filter.SliceFilter
            attribute), 45
slice_stats (ccanalyser.tools.filter.TiledCSliceFilter
            attribute), 52
slice_stats (ccanalyser.tools.filter.TriCSliceFilter
            attribute), 50
slice_stats() (ccanalyser.tools.filter.CCSliceFilter
               property), 48
slice_stats() (ccanalyser.tools.filter.SliceFilter
               property), 46
slice_stats() (ccanal-
                  yser.tools.filter.TiledCSliceFilter
                  attribute), 52
SliceFilter (class in ccanalyser.tools.filter), 45
SLICES
    ccanalyser-alignments-annotate
    command line option, 12
slices (ccanalyser.tools.digest.DigestedRead
        attribute), 44
slices (ccanalyser.tools.filter.CCSliceFilter attribute),
       48
slices (ccanalyser.tools.filter.SliceFilter attribute), 45
slices (ccanalyser.tools.filter.TiledCSliceFilter
        attribute), 51

slices (ccanalyser.tools.filter.TriCSliceFilter
        attribute), 50
SLICES_FN
    ccanalyser-alignments-deduplicate-remove
    command line option, 13
sparse (ccanalyser.tools.pileup.CoolerBedGraph
        attribute), 54
split() (in module ccanalyser.cli.fastq_split), 31
split_intervals_on_chrom() (in module ccanal-
                             yser.utils), 60
statq(ccanalyser.tools.deduplicate.ReadDuplicateRemovalProcess
      attribute), 42
statq (ccanalyser.tools.io.FastqReaderProcess
      attribute), 53
stats_alignment_filtering_collate() (in
                                       module ccanalyser.pipeline.pipeline), 39
stats_deduplication_collate() (in module
                                       ccanalyser.pipeline.pipeline), 38
stats_digestion_collate() (in module ccanal-
                            yser.pipeline.pipeline), 39
stats_trim_collate() (in module ccanal-
                       yser.pipeline.pipeline), 38

T
TiledCSliceFilter (class in ccanal-
                    yser.tools.filter), 51
to_bedtool() (ccanalyser.tools.pileup.CCBedgraph
               method), 55
to_cooler() (ccanalyser.tools.storage.CoolerBinner
               method), 59
to_file() (ccanalyser.tools.pileup.CCBedgraph
               method), 55
to_file() (ccanalyser.tools.pileup.CoolerBedGraph
               method), 55
TriCSliceFilter (class in ccanalyser.tools.filter),
                  50

U
unfiltered (ccanalyser.tools.statistics.DigestionStats
            attribute), 56
UNION_BEDGRAPH
    ccanalyser-reporters-differential
    command line option, 21

Z
zap_files() (in module ccanalyser.utils), 60

```