
MuG - FASTQ Pipelines Documentation

Release 0.1

Mark McDowall

Nov 19, 2018

Contents:

1	Requirements and Installation	1
1.1	Requirements	1
1.2	Installation	2
1.3	Documentation	3
2	Full Installation	5
2.1	Setup the System Environment	5
2.2	Setup pyenv and pyenv-virtualenv	5
2.3	Installation Process	6
2.4	Setup the symlinks	8
2.5	Prepare the Python Environment	11
2.6	Post Installation Tidyup	12
3	Pipelines	13
3.1	Download and index genome files	13
3.2	BioBamBam Alignment Filtering	14
3.3	Bowtie2 Alignment	16
3.4	BSgenome Builder	17
3.5	BS Seeker2 Indexer	18
3.6	BS Seeker2 Aligner	20
3.7	BiSulphate Sequencing Filter	21
3.8	BS Seeker2 Methylation Peak Caller	23
3.9	BWA Alignment - bwa aln	23
3.10	BWA Alignment - bwa mem	24
3.11	ChIP-Seq Analysis	26
3.12	iDamID-Seq Analysis	28
3.13	iNPS	29
3.14	MACS2 Analysis	30
3.15	Mnase-Seq Analysis	32
3.16	RNA-Seq Analysis	33
3.17	TrimGalore	35
3.18	Whole Genome BiSulphate Sequencing Analysis	36
3.19	Hi-C Analysis	36
4	Tools for processing FastQ files	39
4.1	File Validation	39
4.2	Indexers	41

4.3	Aligners	45
4.4	Filters	51
4.5	Peak Calling	54
4.6	Hi-C Parsing	58
5	Utility Functions	71
5.1	Common Functions	71
5.2	Alignment Utilities	71
5.3	Bam Utilities	73
5.4	FASTQ Functions	77
6	Continuous Integration with Travis	81
6.1	Getting Started	81
6.2	Making .travis.yml File	81
6.3	Making harness.sh File	82
6.4	Running Docker container	82
6.5	Setting up Shims	82
7	Setting up and using a Docker Container	83
7.1	Our reason for using a container	83
7.2	Getting Started	83
7.3	Constructing a docker container	84
8	Architectural Design Record	87
8.1	2017-08-15 - Implementation of pigz	87
8.2	2018-01-26 - Disable no-self-use for @tasks	87
8.3	2018-02-28 - BAM Merge Strategy	87
8.4	2018-04-26 - BAM Splitting	88
8.5	2018-05-01 - Compression of FASTQ	88
8.6	2018-05-09 - Handling aligner index decompression	88
8.7	2018-05-22 - GEM Naming	88
8.8	2018-05-22 - TrimGalore	88
8.9	2018-05-31 - Public genomes and indexes	88
8.10	2018-06-01 - Separated WGBS Code Testing	89
8.11	2018-06-01 - Travis Caching	89
8.12	2018-06-04 - Split the WGBS test scripts	89
8.13	2018-06-05 - Use of the logger PROGRESS	89
8.14	2018-06-14 - Paired end alignment	89
8.15	2018-06-18 - Branch tidying during alignment	89
8.16	2018-06-27 - Remove reads marked as duplicate by BioBamBam	90
8.17	2018-07-11 - Changes FASTQ splitter file management	90
8.18	2018-07-16 - Modified handling of file locations	90
8.19	2018-08-02 - Added in Paired End BAM file handling for MACS2	90
8.20	2018-07-16 - Modified handling of file locations	90
8.21	2018-08-07 - Storing tool parameters as part of the metadata	90
8.22	2018-08-07 - Extra output files from MACS2	90
8.23	2018-08-13 - Normalised the use of OSError	91
8.24	2018-08-15 - Use the config.json execution path	91
8.25	2018-08-16 - Prevent further duplicate filtering by MACS2	91
8.26	2018-09-04 - Adding functionality to bam_utils and MACS2	91
8.27	2018-09-17 - Updates to tool and pipeline run()	91
8.28	2018-08-22 - Improvement of tadbit tools wrappers	91
8.29	2018-09-25 - Converting the Kallisto TSV file to BED	91
8.30	2018-10-18 - Multi File handling for the DamID-seq Pipeline	92
8.31	2018-10-25 - WGBS Pipeline Create BigWig files as standard	92

8.32	2018-10-31 - Modify the file names for docker script	92
8.33	2018-11-08 - Modifications for the movement of files	92
8.34	2018-11-16 - Reading of Gzipped FASTQ files	92
9	License	93
10	Testing	97
10.1	Sample Data	97
10.2	Pipelines	109
10.3	Tools	113
11	Indices and tables	117
	Python Module Index	119

Source (github)

1.1 Requirements

1.1.1 Software

- Python 2.7.10+ or 3.6+
- bedtools (specifically bamtobed)
- bedToBigBed - http://hgdownload.cse.ucsc.edu/admin/exe/linux.x86_64/
- wigToBigWig - http://hgdownload.cse.ucsc.edu/admin/exe/linux.x86_64/
- BioBamBam2
- Bowtie2
- BWA
- FastQC
- GEMtools
- HDF5
- iNPS
- Kallisto
- libmaus2
- pyenv
- R 2.9.1+
- SAMtools

- MCL
- pigz

1.1.2 Python Modules

- numpy
- h5py
- pysam
- pyBigWig
- scipy
- matplotlib
- rpy2

Instructions for the following modules are listed in the installation section. All other python modules should be installed with pip prior to the following libraries.

- BS-Seeker2
- TADbit

1.2 Installation

For a guide to the full installation procedure the [Full Installation](#).

Directly from GitHub:

```
1 cd ${HOME}/code
2
3 git clone https://github.com/Multiscale-Genomics/mg-process-fastq.git
4
5 cd mg-process-fastq
```

Create the Python environment

```
1 pyenv-virtualenv 2.7.10 mg-process-fastq
2 pip install --editable .
```

Install the pyTADbit modules

```
1 cd ${HOME}/lib
2 wget https://github.com/3DGenomes/tadbit/archive/master.zip -O tadbit.zip
3 unzip tadbit.zip
4 cd tadbit-master
5
6 pyenv activate mg-process-fastq
7 pip install .
```

Check out the following software for use by the process_wgbs.py pipeline:


```
1 cd cd ${HOME}/lib
2 gti clone https://github.com/BSSeeker/BSseeker2.git
3
4 cd ${HOME}/code
5 cd mg-process-fastq
6 ln -s $code_root/bs_align bs_align
7 ln -s $code_root/bs_index bs_index
8 ln -s $code_root/bs_utils bs_utils
9
10 cd cd ${HOME}/code/mg-process-fastq/tool
11 ln -s $code_root/FilterReads.py FilterReads.py
```

1.3 Documentation

To build the documentation:

```
1 pip install Sphinx
2 pip install sphinx-autobuild
3 cd docs
4 make html
```


The following document is for the full installation of all software required by the mg-process-fastq module and all programmes that it uses. The document has been written with Ubuntu Linux in mind, although many of the commands are cross platform (*nix) compliant.

If you already have certain packages installed feel free to skip over certain steps. Likewise the bin, lib and code directories are relative to the home dir; if this is not the case for your system then make the required changes when running these commands.

2.1 Setup the System Environment

```
1 sudo apt-get install -y make build-essential libssl-dev zlib1g-dev      \\
2 libbz2-dev libreadline-dev libsqlite3-dev wget curl llvm libncurses5-dev \\
3 libncursesw5-dev xz-utils tk-dev unzip mcl libgtk2.0-dev r-base-core    \\
4 libcurl4-gnutls-dev python-rpy2 git libtbb2 pigz liblzma-dev libhdf5-dev \\
5 texlive-latex-base tree libblas-dev liblapack-dev
6
7 cd ${HOME}
8 mkdir bin lib code
9 echo 'export PATH="${HOME}/bin:$PATH"' >> ~/.bash_profile
```

2.2 Setup pyenv and pyenv-virtualenv

This is required for managing the version of Python and the installation environment for the Python modules so that they can be installed in the user space.

```
1 git clone https://github.com/pyenv/pyenv.git ~/.pyenv
2 echo 'export PYENV_ROOT="${HOME}/.pyenv"' >> ~/.bash_profile
3 echo 'export PATH="${PYENV_ROOT}/bin:$PATH"' >> ~/.bash_profile
4 echo 'eval "$(pyenv init -)"' >> ~/.bash_profile
```

(continues on next page)

(continued from previous page)

```
5
6 # Add the .bash_profile to your .bashrc file
7 echo 'source ~/.bash_profile' >> ~/.bashrc
8
9 git clone https://github.com/pyenv/pyenv-virtualenv.git ${PYENV_ROOT}/plugins/pyenv-
  ↳virtualenv
10
11 pyenv install 2.7.12
12 pyenv virtualenv 2.7.12 mg-process-fastq
13
14 # Python 3 environment required by iNPS
15 pyenv install 3.5.3
16 ln -s ${HOME}/.pyenv/versions/3.5.3/bin/python ${HOME}/bin/py3
```

2.3 Installation Process

2.3.1 UCSC Tools

```
1 cd ${HOME}/lib
2 wget http://hgdownload.cse.ucsc.edu/admin/exe/linux.x86_64/bedToBigBed
3 wget http://hgdownload.cse.ucsc.edu/admin/exe/linux.x86_64/wigToBigWig
4
5 wget http://hgdownload.soe.ucsc.edu/admin/exe/linux.x86_64/faToTwoBit
6 wget http://hgdownload.soe.ucsc.edu/admin/exe/linux.x86_64/twoBitInfo
7
8 chmod +x bedToBigBed wigToBigWig faToTwoBit twoBitInfo
```

2.3.2 BioBamBam2

BioBamBam is used for the filtering of aligned reads as part of the ChIP-seq pipeline. It also requires the libmaus2 package to be installed.

```
1 cd ${HOME}/lib
2 git clone https://github.com/gt1/libmaus2.git
3 cd libmaus2
4 libtoolize
5 aclocal
6 autoheader
7 automake --force-missing --add-missing
8 autoconf
9 ./configure --prefix=${HOME}/lib/libmaus2
10 make
11 make install
12
13 cd ${HOME}/lib
14 git clone https://github.com/gt1/biobambam2.git
15 cd biobambam2
16 autoreconf -i -f
17 ./configure --with-libmaus2=${HOME}/lib/libmaus2 --prefix=${HOME}/lib/biobambam2
18 make install
```

2.3.3 Bowtie2 Aligner

```
1 cd ${HOME}/lib
2 wget --max-redirect 1 https://downloads.sourceforge.net/project/bowtie-bio/bowtie2/2.
  ↳3.4/bowtie2-2.3.4-linux-x86_64.zip
3 unzip bowtie2-2.3.4-linux-x86_64.zip
```

2.3.4 BWA Sequence Aligner

```
1 cd ${HOME}/lib
2 git clone https://github.com/lh3/bwa.git
3 cd bwa
4 make
```

2.3.5 FastQC

```
1 cd ${HOME}/lib
2 wget http://www.bioinformatics.babraham.ac.uk/projects/fastqc/fastqc_v0.11.5.zip
3 unzip fastqc_v0.11.5.zip
4 cd FastQC/
5 chmod 755 fastqc
```

2.3.6 GEM Sequence Aligner

```
1 cd ${HOME}/lib
2 wget http://barnaserver.com/gemtools/releases/GEMTools-static-core2-1.7.1.tar.gz
3 tar -xzf GEMTools-static-core2-1.7.1.tar.gz
```

2.3.7 iNPS Peak Caller

```
1 cd ${HOME}/lib
2 mkdir iNPS
3 cd iNPS
4 wget http://www.picb.ac.cn/hanlab/files/iNPS_V1.2.2.zip
5 unzip iNPS_V1.2.2.zip
6
7 cd ${HOME}/bin
8 touch iNPS
9 cat iNPS <<EOL
10 #!/usr/bin/env bash
11 py3 ${HOME}/lib/iNPS/iNPS_V1.2.2.py "$@"
12 EOL
13
14 chmod 777 iNPS
```

2.3.8 Kallisto

```
1 cd ${HOME}/lib
2 wget https://github.com/pachterlab/kallisto/releases/download/v0.43.1/kallisto_linux-
  ↳v0.43.1.tar.gz
3 tar -xzf kallisto_linux-v0.43.1.tar.gz
```

2.3.9 SAMtools

```
1 cd ${HOME}/lib
2 git clone https://github.com/samtools/htslib.git
3 cd htslib
4 autoheader
5 autoconf
6 ./configure --prefix=${HOME}/lib/htslib
7 make
8 make install
9
10 cd ${HOME}/lib
11 git clone https://github.com/samtools/samtools.git
12 cd samtools
13 autoheader
14 autoconf -Wno-syntax
15 ./configure --prefix=${HOME}/lib/samtools
16 make
17 make install
```

2.3.10 bedTools

```
1 cd ${HOME}/lib
2 wget https://github.com/arq5x/bedtools2/releases/download/v2.26.0/bedtools-2.26.0.tar.
  ↳gz
3 tar -zxvf bedtools-2.26.0.tar.gz
4 cd bedtools2
5 make
```

2.4 Setup the symlinks

```
1 cd ${HOME}/bin
2
3 ln -s ${HOME}/lib/bedtools2/bin/bedtools bedtools
4
5 ln -s ${HOME}/lib/bedToBigBed bedToBigBed
6 ln -s ${HOME}/lib/wigToBigWig wigToBigWig
7 ln -s ${HOME}/lib/faToTwoBit faToTwoBit
8 ln -s ${HOME}/lib/twoBitInfo twoBitInfo
9
10 ln -s ${HOME}/lib/bwa/bwa bwa
11
12 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2 bowtie2
```

(continues on next page)

(continued from previous page)

```

13 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-align-s bowtie2-align-s
14 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-align-l bowtie2-align-l
15 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-build bowtie2-build
16 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-build-s bowtie2-build-s
17 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-build-l bowtie2-build-l
18 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-inspect bowtie2-inspect
19 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-inspect-s bowtie2-inspect-s
20 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-inspect-l bowtie2-inspect-l
21
22 ln -s ${HOME}/lib/FastQC/fastqc
23
24 ln -s ${HOME}/lib/gemtools-1.7.1-core2/bin/gem-2-bed gem-2-bed
25 ln -s ${HOME}/lib/gemtools-1.7.1-core2/bin/gem-2-gem gem-2-gem
26 ln -s ${HOME}/lib/gemtools-1.7.1-core2/bin/gem-2-sam gem-2-sam
27 ln -s ${HOME}/lib/gemtools-1.7.1-core2/bin/gem-2-wig gem-2-wig
28 ln -s ${HOME}/lib/gemtools-1.7.1-core2/bin/gem-indexer gem-indexer
29 ln -s ${HOME}/lib/gemtools-1.7.1-core2/bin/gem-indexer_bwt-dna gem-indexer_bwt-dna
30 ln -s ${HOME}/lib/gemtools-1.7.1-core2/bin/gem-indexer_fasta2meta+cont gem-indexer_
  ↪ fasta2meta+cont
31 ln -s ${HOME}/lib/gemtools-1.7.1-core2/bin/gem-indexer_generate gem-indexer_generate
32 ln -s ${HOME}/lib/gemtools-1.7.1-core2/bin/gem-info gem-info
33 ln -s ${HOME}/lib/gemtools-1.7.1-core2/bin/gem-mapper gem-mapper
34 ln -s ${HOME}/lib/gemtools-1.7.1-core2/bin/gemtools gemtools
35
36 ln -s ${HOME}/lib/iNPS/iNPS_V1.2.2.py iNPS_V1.2.2.py
37 ln -s ${HOME}/lib/kallisto_linux-v0.43.1/kallisto kallisto
38
39 ln -s ${HOME}/lib/htslib/bin/bgzip bgzip
40 ln -s ${HOME}/lib/htslib/bin/htsfile htsfile
41 ln -s ${HOME}/lib/htslib/bin/tabix tabix
42
43 ln -s ${HOME}/lib/samtools/bin/ace2sam ace2sam
44 ln -s ${HOME}/lib/samtools/bin/blast2sam.pl blast2sam.pl
45 ln -s ${HOME}/lib/samtools/bin/bowtie2sam.pl bowtie2sam.pl
46 ln -s ${HOME}/lib/samtools/bin/export2sam.pl export2sam.pl
47 ln -s ${HOME}/lib/samtools/bin/interpolate_sam.pl interpolate_sam.pl
48 ln -s ${HOME}/lib/samtools/bin/maq2sam-long maq2sam-long
49 ln -s ${HOME}/lib/samtools/bin/maq2sam-short maq2sam-short
50 ln -s ${HOME}/lib/samtools/bin/md5fa md5fa
51 ln -s ${HOME}/lib/samtools/bin/md5sum-lite md5sum-lite
52 ln -s ${HOME}/lib/samtools/bin/novo2sam.pl novo2sam.pl
53 ln -s ${HOME}/lib/samtools/bin/plot-bamstats plot-bamstats
54 ln -s ${HOME}/lib/samtools/bin/psl2sam.pl psl2sam.pl
55 ln -s ${HOME}/lib/samtools/bin/sam2vcf.pl sam2vcf.pl
56 ln -s ${HOME}/lib/samtools/bin/samtools samtools
57 ln -s ${HOME}/lib/samtools/bin/samtools.pl samtools.pl
58 ln -s ${HOME}/lib/samtools/bin/seq_cache_populate.pl seq_cache_populate.pl
59 ln -s ${HOME}/lib/samtools/bin/soap2sam.pl soap2sam.pl
60 ln -s ${HOME}/lib/samtools/bin/varfilter.py varfilter.py
61 ln -s ${HOME}/lib/samtools/bin/wgsim wgsim
62 ln -s ${HOME}/lib/samtools/bin/wgsim_eval.pl wgsim_eval.pl
63 ln -s ${HOME}/lib/samtools/bin/zoom2sam.pl zoom2sam.pl
64
65 ln -s ${HOME}/lib/biobambam2/bin/bam12auxmerge bam12auxmerge
66 ln -s ${HOME}/lib/biobambam2/bin/bam12split bam12split
67 ln -s ${HOME}/lib/biobambam2/bin/bam12strip bam12strip
68 ln -s ${HOME}/lib/biobambam2/bin/bamadapterclip bamadapterclip

```

(continues on next page)

(continued from previous page)

```
69 ln -s ${HOME}/lib/biobambam2/bin/bamadapterfind bamadapterfind
70 ln -s ${HOME}/lib/biobambam2/bin/bamalignfrac bamalignfrac
71 ln -s ${HOME}/lib/biobambam2/bin/bamauxmerge bamauxmerge
72 ln -s ${HOME}/lib/biobambam2/bin/bamauxsort bamauxsort
73 ln -s ${HOME}/lib/biobambam2/bin/bamcat bamcat
74 ln -s ${HOME}/lib/biobambam2/bin/bamchecksort bamchecksort
75 ln -s ${HOME}/lib/biobambam2/bin/bamclipreinsert bamclipreinsert
76 ln -s ${HOME}/lib/biobambam2/bin/bamcollate bamcollate
77 ln -s ${HOME}/lib/biobambam2/bin/bamcollate2 bamcollate2
78 ln -s ${HOME}/lib/biobambam2/bin/bamdownsamplerandom bamdownsamplerandom
79 ln -s ${HOME}/lib/biobambam2/bin/bamexplode bamexplode
80 ln -s ${HOME}/lib/biobambam2/bin/bamfilteraux bamfilteraux
81 ln -s ${HOME}/lib/biobambam2/bin/bamfilterflags bamfilterflags
82 ln -s ${HOME}/lib/biobambam2/bin/bamfilterheader bamfilterheader
83 ln -s ${HOME}/lib/biobambam2/bin/bamfilterheader2 bamfilterheader2
84 ln -s ${HOME}/lib/biobambam2/bin/bamfilterlength bamfilterlength
85 ln -s ${HOME}/lib/biobambam2/bin/bamfiltermc bamfiltermc
86 ln -s ${HOME}/lib/biobambam2/bin/bamfilternames bamfilternames
87 ln -s ${HOME}/lib/biobambam2/bin/bamfilterrg bamfilterrg
88 ln -s ${HOME}/lib/biobambam2/bin/bamfixmateinformation bamfixmateinformation
89 ln -s ${HOME}/lib/biobambam2/bin/bamflagsplit bamflagsplit
90 ln -s ${HOME}/lib/biobambam2/bin/bamheap2 bamheap2
91 ln -s ${HOME}/lib/biobambam2/bin/bamindex bamindex
92 ln -s ${HOME}/lib/biobambam2/bin/bamintervalcomment bamintervalcomment
93 ln -s ${HOME}/lib/biobambam2/bin/bamintervalcommenthist bamintervalcommenthist
94 ln -s ${HOME}/lib/biobambam2/bin/bamlastfilter bamlastfilter
95 ln -s ${HOME}/lib/biobambam2/bin/bammapdist bammapdist
96 ln -s ${HOME}/lib/biobambam2/bin/bammarkduplicates bammarkduplicates
97 ln -s ${HOME}/lib/biobambam2/bin/bammarkduplicates2 bammarkduplicates2
98 ln -s ${HOME}/lib/biobambam2/bin/bammarkduplicatesopt bammarkduplicatesopt
99 ln -s ${HOME}/lib/biobambam2/bin/bammaskflags bammaskflags
100 ln -s ${HOME}/lib/biobambam2/bin/bammdnm bammdnm
101 ln -s ${HOME}/lib/biobambam2/bin/bammerge bammerge
102 ln -s ${HOME}/lib/biobambam2/bin/bamnumericalindex bamnumericalindex
103 ln -s ${HOME}/lib/biobambam2/bin/bamrank bamrank
104 ln -s ${HOME}/lib/biobambam2/bin/bamranksort bamranksort
105 ln -s ${HOME}/lib/biobambam2/bin/bamrecalculatecigar bamrecalculatecigar
106 ln -s ${HOME}/lib/biobambam2/bin/bamrecompress bamrecompress
107 ln -s ${HOME}/lib/biobambam2/bin/bamreset bamreset
108 ln -s ${HOME}/lib/biobambam2/bin/bamscrapcount bamscrapcount
109 ln -s ${HOME}/lib/biobambam2/bin/bamseqchksum bamseqchksum
110 ln -s ${HOME}/lib/biobambam2/bin/bamsormadup bamsormadup
111 ln -s ${HOME}/lib/biobambam2/bin/bamsort bamsort
112 ln -s ${HOME}/lib/biobambam2/bin/bamsplit bamsplit
113 ln -s ${HOME}/lib/biobambam2/bin/bamsplitdiv bamsplitdiv
114 ln -s ${HOME}/lib/biobambam2/bin/bamstreamingmarkduplicates bamstreamingmarkduplicates
115 ln -s ${HOME}/lib/biobambam2/bin/bamtagconversion bamtagconversion
116 ln -s ${HOME}/lib/biobambam2/bin/bamtofastq bamtofastq
117 ln -s ${HOME}/lib/biobambam2/bin/bamvalidate bamvalidate
118 ln -s ${HOME}/lib/biobambam2/bin/bamzztoname bamzztoname
119 ln -s ${HOME}/lib/biobambam2/bin/fastaexplode fastaexplode
120 ln -s ${HOME}/lib/biobambam2/bin/fastqtobam fastqtobam
121 ln -s ${HOME}/lib/biobambam2/bin/fastqtobampar fastqtobampar
122 ln -s ${HOME}/lib/biobambam2/bin/filtersam filtersam
123 ln -s ${HOME}/lib/biobambam2/bin/kmerprob kmerprob
124 ln -s ${HOME}/lib/biobambam2/bin/lasToBAM lasToBAM
125 ln -s ${HOME}/lib/biobambam2/bin/normalisefasta normalisefasta
```


2.5 Prepare the Python Environment

2.5.1 Install APIs and Pipelines

Checkout the code for the DM API and the mg-process-fastq pipelines:

```

1 cd ${HOME}/code
2 pyenv activate mg-process-fastq
3 pip install git+https://github.com/Multiscale-Genomics/mg-dm-api.git
4 pip install git+https://github.com/Multiscale-Genomics/mg-tool-api.git
5
6 git clone https://github.com/Multiscale-Genomics/mg-process-fastq.git
7 cd mg-process-fastq
8 pip install -e .
9 pip install -r requirements.txt

```

2.5.2 Install MACS2

This should get installed as part of the installation in the mg-process-fastq package, if not then it will need to be installed separately.

For Python 2.7:

```

1 cd ${HOME}/code
2 pyenv activate mg-process-fastq
3 pip install MACS2
4
5 ln -s ${HOME}/.pyenv/versions/mg-process-fastq/bin/macs2 ${HOME}/bin/macs2

```

For Python 3.6:

```

1 cd ${HOME}/code
2 pyenv activate mg-process-fastq
3 git clone https://github.com/taoliu/MACS.git
4 cd MACS
5 git checkout MACS2p3
6
7 cython MACS2/*.pyx
8 cython MACS2/IO/*.pyx
9 python setup_w_cython.py install
10
11 pip install .
12 alias macs2="macs2p3"

```

2.5.3 Install IDEAR

```

1 cd ${HOME}/lib
2 source("https://bioconductor.org/biocLite.R")
3 biocLite("BSgenome")
4 biocLite("DESeq2")
5 if(!require("devtools")) install.packages("devtools")
6 devtools::install_bitbucket("juanlmateo/idear")

```

2.5.4 Install TADbit

```
1 cd ${HOME}/lib
2 wget https://github.com/3DGenomes/TADbit/archive/dev.zip -O tadbit.zip
3 unzip tadbit.zip
4 cd TADbit-dev
5
6 # If the pyenv env is not called mg-process-fastq then change this to match,
7 # the same is true for the version of python
8 python setup.py install --install-lib=${HOME}/.pyenv/versions/mg-process-fastq/lib/
  ↳python2.7/site-packages/ --install-scripts=${HOME}/bin
```

2.5.5 Install BSseeker

```
1 cd ${HOME}/lib
2 git clone https://github.com/BSseeker/BSseeker2.git
3
4 cd ${HOME}/code/mg-process-fastq
5 ln -s ${HOME}/lib/BSseeker2/bs_align bs_align
6 ln -s ${HOME}/lib/BSseeker2/bs_index bs_index
7 ln -s ${HOME}/lib/BSseeker2/bs_utils bs_utils
```

2.5.6 Trim Galore

```
1 cd ${HOME}/lib
2 pip install cutadapt
3 wget -O trim_galore.tar.gz https://github.com/FelixKrueger/TrimGalore/archive/0.5.0.
  ↳tar.gz
4 tar -xzf trim_galore.tar.gz
5
6 cd ${HOME}/bin
7 ln -s ${HOME}/lib/TrimGalore-0.5.0/trim_galore trim_galore
```

Running on a COMPSs VM the symlink will need to be created in a system accessible area:

```
1 sudo ln -s ${HOME}/lib/TrimGalore-0.4.3/trim_galore /usr/local/bin/trim_galore
2 pip install cutadapt
```

2.6 Post Installation Tidyup

```
1 cd ${HOME}/lib
2 rm *.zip *.tar.gz
```

3.1 Download and index genome files

This pipeline is for the indexing of genomes once they have been loaded into the VRE. It indexes each new genome with Bowtie2, BWA and GEM. These indexes can then be used by the other pipelines.

3.1.1 Running from the command line

Parameters

taxon_id [int] Species taxonomic ID

assembly [str] Genomic assembly ID

genome [str] Location of the genomes FASTA file

Returns

Bowtie2 index files BWA index files GEM index file

Example

When running the pipeline on a local machine without COMPSs:

```
1 python process_genome.py \
2   --config tests/json/config_genome_indexer.json \
3   --in_metadata tests/json/input_genome_indexer.json \
4   --out_metadata tests/json/output_genome_indexer.json \
5   --local
```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```
1 runcomps      \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_genome.py \
7     --config tests/json/config_genome_indexer.json \
8     --in_metadata tests/json/input_genome_indexer.json \
9     --out_metadata tests/json/output_genome_indexer.json
```

3.1.2 Methods

class `process_genome.process_genome` (*configuration=None*)

Workflow to download and pre-index a given genome

run (*input_files, metadata, output_files*)

Main run function for the indexing of genome assembly FASTA files. The pipeline uses Bowtie2, BWA and GEM ready for use in pipelines that rely on alignment.

Parameters

- **input_files** (*dict*) –
genome [str] List of file locations
- **metadata** (*dict*) –
genome [dict] Required meta data
- **output_files** (*dict*) –
bwa_index [str] Location of the BWA index archive files
bwt_index [str] Location of the Bowtie2 index archive file
gem_index [str] Location of the GEM index file
genome_gem [str] Location of a the FASTA file generated for the GEM indexing step

Returns

- **outputfiles** (*dict*) – List of locations for the output index files
- **output_metadata** (*dict*) – Metadata about each of the files

3.2 BioBamBam Alignment Filtering

This pipeline to filter sequencing artifacts from aligned reads.

3.2.1 Running from the command line

Parameters

config [str] Configuration JSON file

in_metadata [str] Location of input JSON metadata for files

out_metadata [str] Location of output JSON metadata for files

Returns

filtered [file] Filtered bam file

Example

REQUIREMENT - Needs an aligned bam file

When running the pipeline on a local machine without COMPSs:

```

1 python process_biobambam.py \
2   --config tests/json/config_biobambam.json \
3   --in_metadata tests/json/input_biobambam.json \
4   --out_metadata tests/json/output_biobambam.json \
5   --local

```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```

1 runcompsss \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_biobambam.py \
7   --config tests/json/config_biobambam.json \
8   --in_metadata tests/json/input_biobambam.json \
9   --out_metadata tests/json/output_biobambam.json

```

3.2.2 Methods

class `process_biobambam.process_biobambam` (*configuration=None*)

Functions for filtering FastQ alignments with BioBamBam.

run (*input_files, metadata, output_files*)

Main run function for filtering FastQ aligned reads using BioBamBam.

Parameters

- **input_files** (*dict*) – Location of the initial input files required by the workflow
 - bam** [str] Location of BAM file
- **metadata** (*dict*) – Input file meta data associated with their roles
 - bam** : str
- **output_files** (*dict*) – Output file locations
 - filtered** : str

Returns

- **output_files** (*dict*) – Output file locations associated with their roles, for the output
 - filtered** [str] Filtered version of the bam file
- **output_metadata** (*dict*) – Output metadata for the associated files in `output_files`
 - filtered** : Metadata

3.3 Bowtie2 Alignment

This pipeline aligns FASTQ reads to a given indexed genome. The pipeline can handle single-end and paired-end reads.

3.3.1 Running from the command line

Parameters

config [str] Configuration JSON file

in_metadata [str] Location of input JSON metadata for files

out_metadata [str] Location of output JSON metadata for files

Returns

bam [file] Aligned reads in bam file

Example

REQUIREMENT - Needs the indexing step to be run first

When running the pipeline on a local machine without COMPSs:

```
1 python process_align_bowtie.py \  
2   --config tests/json/config_bowtie2.json \  
3   --in_metadata tests/json/input_bowtie2.json \  
4   --out_metadata tests/json/output_bowtie2.json \  
5   --local
```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```
1 runcompss \  
2   --lang=python \  
3   --library_path=${HOME}/bin \  
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \  
5   --log_level=debug \  
6   process_align_bowtie.py \  
7     --config tests/json/config_bowtie2_single.json \  
8     --in_metadata tests/json/input_bowtie2_single_metadata.json \  
9     --out_metadata tests/json/output_bowtie2_single.json
```

```
1 runcompss \  
2   --lang=python \  
3   --library_path=${HOME}/bin \  
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \  
5   --log_level=debug \  
6   process_align_bowtie.py \  
7     --config tests/json/config_bowtie2_paired.json \  
8     --in_metadata tests/json/input_bowtie2_paired_metadata.json \  
9     --out_metadata tests/json/output_bowtie2_paired.json
```

3.3.2 Methods

class `process_align_bowtie.process_bowtie` (*configuration=None*)
 Functions for aligning FastQ files with Bowtie2

run (*input_files, metadata, output_files*)
 Main run function for aligning FastQ reads with Bowtie2.

Currently this can only handle a single data file and a single background file.

Parameters

- **input_files** (*dict*) – Location of the initial input files required by the workflow
 - genome** [str] Genome FASTA file
 - index** [str] Location of the BWA archived index files
 - loc** [str] Location of the FASTQ reads files
 - fastq2** [str] [OPTIONAL] Location of the FASTQ reads file for paired end data
- **metadata** (*dict*) – Input file meta data associated with their roles
 - genome** : str **index** : str **loc** : str **fastq2** : str
- **output_files** (*dict*) – Output file locations
 - bam** [str] Output bam file location

Returns

- **output_files** (*dict*) – Output file locations associated with their roles, for the output
 - bam** [str] Aligned FASTQ short read file locations
- **output_metadata** (*dict*) – Output metadata for the associated files in `output_files`
 - bam** : Metadata

3.4 BSgenome Builder

This pipeline can process FASTQ to identify protein-DNA binding sites.

3.4.1 Running from the command line

Parameters

config [str] Configuration JSON file
in_metadata [str] Location of input JSON metadata for files
out_metadata [str] Location of output JSON metadata for files

Returns

bsgenome [file] BSgenome index
genome_2bit [file] Compressed representation of the genome required for generating the index
chrom_size [file] Location of the chrom.size file

seed_file [file] Configuration file for generating the BSgenome R package

Example

When running the pipeline on a local machine without COMPSs:

```
1 python process_bsgenome.py \
2   --config tests/json/config_bsgenome.json \
3   --in_metadata tests/json/input_bsgenome.json \
4   --out_metadata tests/json/output_bsgenome.json \
5   --local
```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```
1 runcompss \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_bsgenome.py \
7   --config tests/json/config_bsgenome.json \
8   --in_metadata tests/json/input_bsgenome.json \
9   --out_metadata tests/json/output_bsgenome.json
```

3.4.2 Methods

class `process_bsgenome.process_bsgenome` (*configuration=None*)
Workflow to download and pre-index a given genome

run (*input_files, metadata, output_files*)

Main run function for the indexing of genome assembly FASTA files. The pipeline uses Bowtie2, BWA and GEM ready for use in pipelines that rely on alignment.

Parameters

- **input_files** (*dict*) –
genome [str] Location of the FASTA input file
- **metadata** (*dict*) –
genome [dict] Required meta data
- **output_files** (*dict*) –
BSgenome [str] Location of a the BSgenome R package

Returns

- **outputfiles** (*dict*) – List of locations for the output index files
- **output_metadata** (*dict*) – Metadata about each of the files

3.5 BS Seeker2 Indexer

This pipeline can process FASTQ to identify protein-DNA binding sites.

3.5.1 Running from the command line

Parameters

config [str] Configuration JSON file

in_metadata [str] Location of input JSON metadata for files

out_metadata [str] Location of output JSON metadata for files

Returns

index [file] BS Seeker2 index

Example

When running the pipeline on a local machine without COMPSs:

```

1 python process_bs_seeker_index.py \
2   --config tests/json/config_wgbs_index.json \
3   --in_metadata tests/json/input_wgbs_index_metadata.json \
4   --out_metadata tests/json/output_wgbs_index.json \
5   --local

```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```

1 runcomps \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_bs_seeker_index.py \
7   --config tests/json/config_wgbs_index.json \
8   --in_metadata tests/json/input_wgbs_index_metadata.json \
9   --out_metadata tests/json/output_wgbs_index.json

```

3.5.2 Methods

class `process_bs_seeker_index.process_bs_seeker_index` (*configuration=None*)
 Functions for aligning FastQ files with BWA

run (*input_files, metadata, output_files*)

Main run function for generatigng the index files required by BS Seeker2.

Parameters

- **input_files** (*dict*) – List of strings for the locations of files. These should include:
 - genome_fa** [str] Genome assembly in FASTA
- **metadata** (*dict*) – Input file meta data associated with their roles
 - genome** : str
- **output_files** (*dict*) – Output file locations
 - bam** [str] Output bam file location

Returns

output_files – Output file locations associated with their roles, for the output
index : str

Return type dict

3.6 BS Seeker2 Aligner

This pipeline aligns FASTQ paired end reads using BS Seeker2 and Bowtie2.

3.6.1 Running from the command line

Parameters

config [str] Configuration JSON file

in_metadata [str] Location of input JSON metadata for files

out_metadata [str] Location of output JSON metadata for files

Returns

bam [file] Aligned Bam file

bai [file] Aligned Bam index file

Example

When running the pipeline on a local machine without COMPSs:

```
1 python process_bs_seeker_aligner.py \
2   --config tests/json/config_wgbs_align.json \
3   --in_metadata tests/json/input_wgbs_align_metadata.json \
4   --out_metadata tests/json/output_wgbs_align.json \
5   --local
```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```
1 runcompss \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_bs_seeker_aligner.py \
7   --config tests/json/config_wgbs_align.json \
8   --in_metadata tests/json/input_wgbs_align_metadata.json \
9   --out_metadata tests/json/output_wgbs_align.json
```

3.6.2 Methods

class `process_bs_seeker_aligner.process_bs_seeker_aligner` (*configuration=None*)
 Functions for downloading and processing whole genome bisulfate sequencings (WGBS) files. Files are filtered, aligned and analysed for points of methylation

run (*input_files, metadata, output_files*)

This pipeline processes paired-end FASTQ files to identify methylated regions within the genome.

Parameters

- **input_files** (*dict*) – List of strings for the locations of files. These should include:
 - genome_fa** [str] Genome assembly in FASTA
 - fastq1** [str] Location for the first filtered FASTQ file for single or paired end reads
 - fastq2** [str] Location for the second filtered FASTQ file if paired end reads
 - index** [str] Location of the index file
- **metadata** (*dict*) – Input file meta data associated with their roles
 - `genome_fa` : dict `fastq1` : dict `fastq2` : dict `index` : dict
- **output_files** (*dict*) – bam : str bai : str

Returns `bam/bai` – Location of the alignment bam file and the associated index

Return type str

3.7 BiSulphate Sequencing Filter

This pipeline processes FASTQ files to filter out duplicate reads.

3.7.1 Running from the command line

Parameters

config [str] Configuration JSON file

in_metadata [str] Location of input JSON metadata for files

out_metadata [str] Location of output JSON metadata for files

Returns

fastq1_filtered/fastq1_filtered [str] Locations of the filtered FASTQ files from which alignments were made

fastq2_filtered/fastq2_filtered [str] Locations of the filtered FASTQ files from which alignments were made

Example

When running the pipeline on a local machine without COMPSs:

```
1 python process_bs_seeker_filter.py \
2   --config tests/json/config_wgbs_filter.json \
3   --in_metadata tests/json/input_wgbs_filter_metadata.json \
4   --out_metadata tests/json/output_metadata.json \
5   --local
```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```
1 runcompss \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_bs_seeker_filter.py \
7   --config tests/json/config_wgbs_filter.json \
8   --in_metadata tests/json/input_wgbs_filter_metadata.json \
9   --out_metadata tests/json/output_metadata.json
```

3.7.2 Methods

class `process_bs_seeker_filter.process_bsFilter` (*configuration=None*)

Functions for filtering FASTQ files. Files are filtered for removal of duplicate reads. Low quality reads in qseq file can also be filtered.

run (*input_files, metadata, output_files*)

This pipeline processes FASTQ files to filter duplicate entries

Parameters

- **input_files** (*dict*) – List of strings for the locations of files. These should include:
 - fastq1** [str] Location for the first FASTQ file for single or paired end reads
 - fastq2** [str] Location for the second FASTQ file if paired end reads [OPTIONAL]
- **metadata** (*dict*) – Input file meta data associated with their roles
 - fastq1** : str
 - fastq2** [str] [OPTIONAL]
- **output_files** (*dict*) – **fastq1_filtered** : str
 - fastq2_filtered** [str] [OPTIONAL]

Returns

- **fastq1_filtered/fastq1_filtered** (*str*) – Locations of the filtered FASTQ files from which alignments were made
- **fastq2_filtered/fastq2_filtered** (*str*) – Locations of the filtered FASTQ files from which alignments were made

3.8 BS Seeker2 Methylation Peak Caller

3.9 BWA Alignment - bwa aln

This pipeline aligns FASTQ reads to a given indexed genome. The pipeline can handle single-end and paired-end reads.

3.9.1 Running from the command line

Parameters

config [str] Configuration JSON file

in_metadata [str] Location of input JSON metadata for files

out_metadata [str] Location of output JSON metadata for files

Returns

bam [file] Aligned reads in bam file

Example

REQUIREMENT - Needs the indexing step to be run first

When running the pipeline on a local machine without COMPSs:

```

1 python process_align_bwa.py \
2   --config tests/json/config_chipseq.json \
3   --in_metadata tests/json/input_chipseq.json \
4   --out_metadata tests/json/output_chipseq.json \
5   --local

```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```

1 runcompss \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_align_bwa.py \
7   --config tests/json/config_bwa_aln_single.json \
8   --in_metadata tests/json/input_bwa_aln_single_metadata.json \
9   --out_metadata tests/json/output_bwa_aln_single.json

```

```

1 runcompss \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_align_bwa.py \
7   --config tests/json/config_bwa_aln_paired.json \

```

(continues on next page)

(continued from previous page)

```
8 --in_metadata tests/json/input_bwa_aln_paired_metadata.json \  
9 --out_metadata tests/json/output_bwa_aln_paired.json
```

3.9.2 Methods

class `process_align_bwa.process_bwa` (*configuration=None*)
Functions for aligning FastQ files with BWA ALN

run (*input_files, metadata, output_files*)
Main run function for aligning FastQ reads with BWA ALN.

Parameters

- **input_files** (*dict*) – Location of the initial input files required by the workflow
 - genome** [str] Genome FASTA file
 - index** [str] Location of the BWA archived index files
 - loc** [str] Location of the FASTQ reads files
 - fastq2** [str] [OPTIONAL] Location of the FASTQ reads file for paired end data
- **metadata** (*dict*) – Input file meta data associated with their roles
 - genome : str index : str loc : str fastq2 : str
- **output_files** (*dict*) – Output file locations
 - bam** [str] Output bam file location

Returns

- **output_files** (*dict*) – Output file locations associated with their roles, for the output
 - bam** [str] Aligned FASTQ short read file locations
- **output_metadata** (*dict*) – Output metadata for the associated files in `output_files`
 - bam : Metadata

3.10 BWA Alignment - bwa mem

This pipeline aligns FASTQ reads to a given indexed genome. The pipeline can handle single-end and paired-end reads.

3.10.1 Running from the command line

Parameters

config [str] Configuration JSON file
in_metadata [str] Location of input JSON metadata for files
out_metadata [str] Location of output JSON metadata for files

Returns

bam [file] Aligned reads in bam file

Example

REQUIREMENT - Needs the indexing step to be run first

When running the pipeline on a local machine without COMPSs:

```

1 python process_align_bwa.py \
2   --config tests/json/config_chipseq.json \
3   --in_metadata tests/json/input_chipseq.json \
4   --out_metadata tests/json/output_chipseq.json \
5   --local

```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```

1 runcomps \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_align_bwa_mem.py \
7   --config tests/json/config_bwa_mem_single.json \
8   --in_metadata tests/json/input_bwa_mem_single_metadata.json \
9   --out_metadata tests/json/output_bwa_mem_single.json

```

```

1 runcomps \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_align_bwa_mem.py \
7   --config tests/json/config_bwa_mem_paired.json \
8   --in_metadata tests/json/input_bwa_mem_paired_metadata.json \
9   --out_metadata tests/json/output_bwa_mem_paired.json

```

3.10.2 Methods

class `process_align_bwa_mem.process_bwa_mem` (*configuration=None*)

Functions for aligning FastQ files with BWA MEM

run (*input_files, metadata, output_files*)

Main run function for aligning FastQ data with BWA MEM.

Parameters

- **input_files** (*dict*) – Location of the initial input files required by the workflow
- **genome** [str] Genome FASTA file
- **index** [str] Location of the BWA archived index files
- **loc** [str] Location of the FASTQ reads files
- **fastq2** [str] [OPTIONAL] Location of the FASTQ reads file for paired end data

- **metadata** (*dict*) – Input file meta data associated with their roles
genome : str index : str loc : str fastq2 : str
- **output_files** (*dict*) – Output file locations
bam [str] Output bam file location

Returns

- **output_files** (*dict*) – Output file locations associated with their roles, for the output
bam [str] Aligned FASTQ short read file locations
- **output_metadata** (*dict*) – Output metadata for the associated files in output_files
bam : Metadata

3.11 ChIP-Seq Analysis

This pipeline can process FASTQ to identify protein-DNA binding sites.

3.11.1 Running from the command line

Parameters

- config** [str] Configuration JSON file
- in_metadata** [str] Location of input JSON metadata for files
- out_metadata** [str] Location of output JSON metadata for files

Returns

bed [file] Bed files with the locations of transcription factor binding sites within the genome

Example

REQUIREMENT - Needs the indexing step to be run first

When running the pipeline on a local machine without COMPSs:

```
1 python process_chipseq.py \
2   --config tests/json/config_chipseq.json \
3   --in_metadata tests/json/input_chipseq.json \
4   --out_metadata tests/json/output_chipseq.json \
5   --local
```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```
1 runcompss \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
```

(continues on next page)

(continued from previous page)

```

6 process_chipseq.py      \
7   --config tests/json/config_chipseq.json \
8   --in_metadata tests/json/input_chipseq.json \
9   --out_metadata tests/json/output_chipseq.json

```

3.11.2 Methods

class process_chipseq.**process_chipseq** (*configuration=None*)

Functions for processing Chip-Seq FastQ files. Files are the aligned, filtered and analysed for peak calling

run (*input_files, metadata, output_files*)

Main run function for processing ChIP-seq FastQ data. Pipeline aligns the FASTQ files to the genome using BWA. MACS 2 is then used for peak calling to identify transcription factor binding sites within the genome.

Currently this can only handle a single data file and a single background file.

Parameters

- **input_files** (*dict*) – Location of the initial input files required by the workflow
 - genome** [str] Genome FASTA file
 - index** [str] Location of the BWA archived index files
 - loc** [str] Location of the FASTQ reads files
 - fastq2** [str] Location of the paired end FASTQ file [OPTIONAL]
 - bg_loc** [str] Location of the background FASTQ reads files [OPTIONAL]
 - fastq2_bg** [str] Location of the paired end background FASTQ reads files [OPTIONAL]
- **metadata** (*dict*) – Input file meta data associated with their roles
 - genome : str index : str
 - bg_loc** [str] [OPTIONAL]
- **output_files** (*dict*) – Output file locations
 - bam [, “bam_bg”] : str filtered [, “filtered_bg”] : str narrow_peak : str summits : str
 - broad_peak : str gapped_peak : str

Returns

- **output_files** (*dict*) – Output file locations associated with their roles, for the output
 - bam** [, “bam_bg”] [str] Aligned FASTQ short read file [and aligned background file] locations
 - filtered** [, “filtered_bg”] [str] Filtered versions of the respective bam files
 - narrow_peak** [str] Results files in bed4+1 format
 - summits** [str] Results files in bed6+4 format
 - broad_peak** [str] Results files in bed6+3 format
 - gapped_peak** [str] Results files in bed12+3 format

- **output_metadata** (*dict*) – Output metadata for the associated files in `output_files`
bam [, “bam_bg”] : Metadata filtered [, “filtered_bg”] : Metadata narrow_peak : Metadata
summits : Metadata broad_peak : Metadata gapped_peak : Metadata

3.12 iDamID-Seq Analysis

This pipeline can process FASTQ to identify protein-DNA binding sites.

3.12.1 Running from the command line

Parameters

config [str] Configuration JSON file

in_metadata [str] Location of input JSON metadata for files

out_metadata [str] Location of output JSON metadata for files

Returns

bigwig [file] Bigwig file of the binding profile of transcription factors

Example

REQUIREMENT - Needs the indexing step to be run first

When running the pipeline on a local machine without COMPSs:

```
1 python process_damidseq.py \
2   --config tests/json/config_idamidseq.json \
3   --in_metadata tests/json/input_idamidseq.json \
4   --out_metadata tests/json/output_idamidseq.json \
5   --local
```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```
1 runcomps \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_damidseq.py \
7   --config tests/json/config_idamidseq.json \
8   --in_metadata tests/json/input_idamidseq.json \
9   --out_metadata tests/json/output_idamidseq.json
```

3.12.2 Methods

class `process_damidseq.process_damidseq` (*configuration=None*)

Functions for processing Chip-Seq FastQ files. Files are the aligned, filtered and analysed for peak calling

run (*input_files*, *metadata*, *output_files*)

Main run function for processing DamID-seq FastQ data. Pipeline aligns the FASTQ files to the genome using BWA. iDEAR is then used for peak calling to identify transcription factor binding sites within the genome.

Currently this can only handle a single data file and a single background file.

Parameters

- **input_files** (*dict*) – Location of the initial input files required by the workflow
 - genome** [str] Genome FASTA file
 - index** [str] Location of the BWA archived index files
 - fastq_1** [str] Location of the FASTQ reads files
 - fastq_2** [str] Location of the FASTQ repeat reads files
 - bg_fastq_1** [str] Location of the background FASTQ reads files
 - bg_fastq_2** [str] Location of the background FASTQ repeat reads files
- **metadata** (*dict*) – Input file meta data associated with their roles
 - genome : str index : str fastq_1 : str fastq_2 : str bg_fastq_1 : str bg_fastq_2 : str
- **output_files** (*dict*) – Output file locations
 - bam [, “bam_bg”] : str filtered [, “filtered_bg”] : str

Returns

- **output_files** (*dict*) – Output file locations associated with their roles, for the output
 - bam** [, “bam_bg”] [str] Aligned FASTQ short read file [and aligned background file] locations
 - filtered** [, “filtered_bg”] [str] Filtered versions of the respective bam files
 - bigwig** [str] Location of the bigwig peaks
- **output_metadata** (*dict*) – Output metadata for the associated files in output_files
 - bam [, “bam_bg”] : Metadata filtered [, “filtered_bg”] : Metadata bigwig : Metadata

3.13 iNPS

This pipeline can process bam file to identify nucleosome positions.

3.13.1 Running from the command line

Parameters

config [str] Configuration JSON file

in_metadata [str] Location of input JSON metadata for files

out_metadata [str] Location of output JSON metadata for files

Returns

bed [file] Bed files with the locations of nucleosome binding sites within the genome

Example

REQUIREMENT - Needs the indexing step to be run first

When running the pipeline on a local machine without COMPSs:

```
1 python process_iNPS.py \
2   --config tests/json/config_inps.json \
3   --in_metadata tests/json/input_iNPS_metadata.json \
4   --out_metadata tests/json/output_iNPS.json \
5   --local
```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```
1 runcomps \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_iNPS.py \
7   --config tests/json/config_inps.json \
8   --in_metadata tests/json/input_iNPS_metadata.json \
9   --out_metadata tests/json/output_iNPS.json
```

3.13.2 Methods

class `process_iNPS.process_iNPS` (*configuration=None*)

Functions for improved nucleosome positioning algorithm (iNPS). Bam Files are analysed for peaks for nucleosome positioning

run (*input_files, metadata, output_files*)

This pipeline processes bam files to identify nucleosome regions within the genome and generates bed files.

Parameters

- **input_files** (*dict*) – `bam_file` : str Location of the aligned sequences in bam format
- **output_files** (*dict*) – `peak_bed` : str Location of the collated bed file of nucleosome peak calls

Returns `peak_bed` – Location of the collated bed file of nucleosome peak calls

Return type str

3.14 MACS2 Analysis

Transcript binding site peak caller for ChIP-seq data

3.14.1 Running from the command line

Parameters

config [str] Configuration JSON file

in_metadata [str] Location of input JSON metadata for files

out_metadata [str] Location of output JSON metadata for files

Returns

bam [file] Aligned reads in bam file

Example

REQUIREMENT - Needs the indexing step to be run first

When running the pipeline on a local machine without COMPSs:

```

1 python process_align_bwa.py \
2   --config tests/json/config_mac2.json \
3   --in_metadata tests/json/input_mac2.json \
4   --out_metadata tests/json/output_mac2.json \
5   --local

```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```

1 runcomps \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_mac2.py \
7   --config tests/json/config_mac2_single.json \
8   --in_metadata tests/json/input_mac2_metadata.json \
9   --out_metadata tests/json/output_mac2.json

```

```

1 runcomps \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_mac2.py \
7   --config tests/json/config_mac2_bgd_paired.json \
8   --in_metadata tests/json/input_mac2_bgd_paired_metadata.json \
9   --out_metadata tests/json/output_mac2_bgd.json

```

3.14.2 Methods

class `process_mac2.process_mac2` (*configuration=None*)

Functions for processing Chip-Seq FastQ files. Files are the aligned, filtered and analysed for peak calling

run (*input_files*, *metadata*, *output_files*)

Main run function for processing ChIP-seq FastQ data. Pipeline aligns the FASTQ files to the genome using BWA. MACS 2 is then used for peak calling to identify transcription factor binding sites within the genome.

Currently this can only handle a single data file and a single background file.

Parameters

- **input_files** (*dict*) – Location of the initial input files required by the workflow
 - bam** [str] Location of the aligned reads file
 - bam_bg** [str] Location of the background aligned FASTQ reads file [OPTIONAL]
- **metadata** (*dict*) – Input file meta data associated with their roles
 - bam** : str
 - bam_bg** [str] [OPTIONAL]
- **output_files** (*dict*) – Output file locations
 - narrow_peak** : str **summits** : str **broad_peak** : str **gapped_peak** : str

Returns

- **output_files** (*dict*) – Output file locations associated with their roles, for the output
 - narrow_peak** [str] Results files in bed4+1 format
 - summits** [str] Results files in bed6+4 format
 - broad_peak** [str] Results files in bed6+3 format
 - gapped_peak** [str] Results files in bed12+3 format
- **output_metadata** (*dict*) – Output metadata for the associated files in *output_files*
 - narrow_peak** : Metadata **summits** : Metadata **broad_peak** : Metadata **gapped_peak** : Metadata

3.15 Mnase-Seq Analysis

This pipeline can process FASTQ to identify nucleosome binding sites.

3.15.1 Running from the command line

Parameters

- config** [str] Configuration JSON file
- in_metadata** [str] Location of input JSON metadata for files
- out_metadata** [str] Location of output JSON metadata for files

Returns

- bed** [file] Bed files with the locations of nucleosome binding sites within the genome

Example

REQUIREMENT - Needs the indexing step to be run first

When running the pipeline on a local machine without COMPSs:

```

1 python process_mnaseseq.py \
2   --config tests/json/config_mnaseseq.json \
3   --in_metadata tests/json/input_mnaseseq.json \
4   --out_metadata tests/json/output_mnaseseq.json \
5   --local

```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```

1 runcompss \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_mnaseseq.py \
7   --config tests/json/config_mnaseseq.json \
8   --in_metadata tests/json/input_mnaseseq.json \
9   --out_metadata tests/json/output_mnaseseq.json

```

3.15.2 Methods

class `process_mnaseseq.process_mnaseseq` (*configuration=None*)

Functions for downloading and processing Mnase-seq FastQ files. Files are downloaded from the European Nucleotide Archive (ENA), then aligned, filtered and analysed for peak calling

run (*input_files, metadata, output_files*)

Main run function for processing MNase-Seq FastQ data. Pipeline aligns the FASTQ files to the genome using BWA. iNPS is then used for peak calling to identify nucleosome position sites within the genome.

Parameters

- **files_ids** (*list*) – List of file locations
- **metadata** (*list*) – Required meta data

Returns **outputfiles** – List of locations for the output bam, bed and tsv files

Return type list

3.16 RNA-Seq Analysis

This pipeline can process FASTQ to quantify the level of expression of cDNAs.

3.16.1 Running from the command line

Parameters

config [str] Configuration JSON file

in_metadata [str] Location of input JSON metadata for files

out_metadata [str] Location of output JSON metadata for files

Returns

bed [file] WIG file with the levels of expression for genes

Example

When running the pipeline on a local machine without COMPSs:

```
1 python process_rnaseq.py \
2   --config tests/json/config_rnaseq.json \
3   --in_metadata tests/json/input_rnaseq.json \
4   --out_metadata tests/json/output_rnaseq.json \
5   --local
```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```
1 runcomps \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=/<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_rnaseq.py \
7   --config tests/json/config_rnaseq.json \
8   --in_metadata tests/json/input_rnaseq.json \
9   --out_metadata tests/json/output_rnaseq.json
```

3.16.2 Methods

class `process_rnaseq.process_rnaseq` (*configuration=None*)

Functions for downloading and processing RNA-seq FastQ files. Files are downloaded from the European Nucleotide Archive (ENA), then they are mapped to quantify the amount of cDNA

run (*input_files, metadata, output_files*)

Main run function for processing RNA-Seq FastQ data. Pipeline aligns the FASTQ files to the genome using Kallisto. Kallisto is then also used for peak calling to identify levels of expression.

Parameters

- **files_ids** (*dict*) – List of file locations (genome FASTA, FASTQ_01, FASTQ_02 (for paired ends))
- **metadata** (*list*) – Required meta data
- **output_files** (*list*) – List of output file locations
- **input_files** (*list*) – List of file locations
- **metadata** – Required meta data
- **output_files** – List of output file locations

Returns **outputfiles** – List of locations for the output bam, bed and tsv files

Return type list

Returns

- **outputfiles** (*dict*) – List of locations for the output index files
- **output_metadata** (*dict*) – Metadata about each of the files

3.17 TrimGalore

This pipeline can process FASTQ to trim poor base quality or adapter contamination.

3.17.1 Running from the command line

Parameters

config [str] Configuration JSON file

in_metadata [str] Location of input JSON metadata for files

out_metadata [str] Location of output JSON metadata for files

Returns

fastq_trimmed [file] Location of a fastq file containing the sequences after poor base qualities or contamination trimming

A full description of the Trim Galore files can be found at <https://github.com/FelixKrueger/TrimGalore>

Example

When running the pipeline on a local machine without COMPSs:

```

1 python process_trim_galore.py \
2   --config tests/json/config_trimgalore.json \
3   --in_metadata tests/json/input_trimgalore_metadata.json \
4   --out_metadata tests/json/output_trimgalore.json \
5   --local

```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```

1 runcomps \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_trim_galore.py \
7   --config tests/json/config_trimgalore.json \
8   --in_metadata tests/json/input_trimgalore_metadata.json \
9   --out_metadata tests/json/output_trimgalore.json

```

3.17.2 Methods

class `process_trim_galore.process_trim_galore` (*configuration=None*)

Functions for filtering FASTQ files. Files are filtered for removal of duplicate reads. Low quality reads in qseq file can also be filtered.

run (*input_files, metadata, output_files*)

This pipeline processes FASTQ files to trim low quality base calls and adapter sequences

Parameters `input_files` (*dict*) – List of strings for the locations of files. These should include:

fastq [str] Location for the first FASTQ file for single or paired end reads

metadata [dict] Input file meta data associated with their roles

`output_files` : dict

`fastq_trimmed` : str

Returns `fastq_trimmed`/`fastq_trimmed` – Locations of the filtered FASTQ files from which trimmings were made

Return type str

3.18 Whole Genome BiSulphate Sequencing Analysis

3.19 Hi-C Analysis

This pipeline can process paired end FASTQ files to identify structural interactions that occur so that the genome can fold into the confines of the nucleus

3.19.1 Running from the command line

Parameters

genome [str] Location of the genomes FASTA file

genome_gem [str] Location of the genome GEM index file

taxon_id [int] Species taxonomic ID

assembly [str] Genomic assembly ID

file1 [str] Location of FASTQ file 1

file2 [str] Location of FASTQ file 2

resolutions [str] Comma separated list of resolutions to calculate the matrix for. [DEFAULT : 1000000,10000000]

enzyme_name [str] Name of the enzyme used to digest the genome (example ‘MboI’)

window_type [str] iter | frag. Analysis windowing type to use

windows1 [str] FASTQ sampling window sizes to use for the first paired end FASTQ file, the default is to use `[[1,25], [1,50], [1,75], [1,100]]`. This would be represented as `1,25,50,75,100` as input for this variable

windows2 [str] FASTQ sampling window sizes to use for the second paired end FASTQ file, the default is to use `[[1,25], [1,50], [1,75], [1,100]]`. This would be represented as `1,25,50,75,100` as input for this variable

normalized [int] 1|0. Determines whether the counts of alignments should be normalized

tag [str] Name for the experiment output files to use

Returns

Adjacency List : file
HDF5 Adjacency Array : file

Example

REQUIREMENT - Needs the indexing step to be run first

When running the pipeline on a local machine:

```

1 python process_hic.py \
2   --genome /<dataset_dir>/Homo_sapiens.GRCh38.fasta \
3   --genome_gem /<dataset_dir>/Homo_sapiens.GRCh38.gem \
4   --assembly GCA_000001405.25 \
5   --taxon_id 9606 \
6   --file1 /<dataset_dir>/<file_name>_1.fastq \
7   --file2 /<dataset_dir>/<file_name>_2.fastq \
8   --resolutions 1000000,10000000 \
9   --enzyme_name MboI \
10  --windows1 1,100 \
11  --windows2 1,100 \
12  --normalized 1 \
13  --tag Human.SRR1658573 \
14  --window_type frag

```

When using a local version of the [COMPS virtual machine](<https://www.bsc.es/research-and-development/software-and-apps/software-list/comp-superscalar/>):

```

1 runcomps \
2   --lang=python \
3   --library_path=${HOME}/bin \
4   --pythonpath=/<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
5   --log_level=debug \
6   process_hic.py \
7   --taxon_id 9606 \
8   --genome /<dataset_dir>/Human.GCA_000001405.22_gem.fasta \
9   --assembly GRCh38 \
10  --file1 /<dataset_dir>/Human.SRR1658573_1.fastq \
11  --file2 /<dataset_dir>/Human.SRR1658573_2.fastq \
12  --genome_gem /<dataset_dir>/Human.GCA_000001405.22_gem.fasta.gem \
13  --enzyme_name MboI \
14  --resolutions 10000,100000 \
15  --windows1 1,100 \
16  --windows2 1,100 \
17  --normalized 1 \
18  --tag Human.SRR1658573 \
19  --window_type frag

```

3.19.2 Methods

class `process_hic.process_hic` (*configuration=None*)

Functions for downloading and processing Mnase-seq FastQ files. Files are downloaded from the European Nucleotide Archive (ENA), then aligned, filtered and analysed for peak calling

run (*input_files, metadata, output_files*)

Main run function for processing MNase-Seq FastQ data. Pipeline aligns the FASTQ files to the genome using BWA. iNPS is then used for peak calling to identify nucleosome position sites within the genome.

Parameters

- **files_ids** (*list*) – List of file locations
- **metadata** (*list*) – Required meta data
- **output_files** (*list*) – List of output file locations

Returns **outputfiles** – List of locations for the output bam, bed and tsv files

Return type list

4.1 File Validation

Pipelines and functions assessing the quality of input files.

4.1.1 FastQC

class `tool.validate_fastqc.fastqcTool` (*configuration=None*)
Tool for running indexers over a genome FASTA file

run (*input_files, input_metadata, output_files*)
Tool for assessing the quality of reads in a FastQ file

Parameters

- **input_files** (*dict*) –
fastq [str] List of file locations
- **metadata** (*dict*) –
fastq [dict] Required meta data
- **output_files** (*dict*) –
report [str] Location of the HTML

Returns **array** – First element is a list of the index files. Second element is a list of the matching metadata

Return type list

validate (***kwargs*)
FastQC Validator

Parameters

- **FastQC_file** (*str*) – Location of the FastQ file

- **report_loc** (*str*) – Location of the output report file

4.1.2 TrimGalore

class `tool.trimgalore.trimgalore` (*configuration=None*)

Tool for trimming FASTQ reads that are of low quality

static `get_trimgalore_params` (*params*)

Function to handle for extraction of commandline parameters

Parameters `params` (*dict*) –

Returns

Return type list

run (*input_files, input_metadata, output_files*)

The main function to run TrimGalore to remove low quality and very short reads. TrimGalore uses CutAdapt and FASTQC for the analysis.

Parameters

- **input_files** (*dict*) –
 - fastq1** [string] Location of the FASTQ file
 - fastq2** [string] [OPTIONAL] Location of the paired end FASTQ file
- **metadata** (*dict*) – Matching metadata for the input FASTQ files

Returns

- **output_files** (*dict*) –
 - fastq1_trimmed** [str] Location of the trimmed FASTQ file
 - fastq2_trimmed** [str] [OPTIONAL] Location of a trimmed paired end FASTQ file
- **output_metadata** (*dict*) – Matching metadata for the output files

trimgalore_paired (***kwargs*)

Trims and removes low quality subsections and reads from paired-end FASTQ files

Parameters

- **fastq_file_in** (*str*) – Location of the input fastq file
- **fastq_file_out** (*str*) – Location of the output fastq file
- **params** (*dict*) – Parameters to use in TrimGalore

Returns Indicator of the success of the function

Return type bool

trimgalore_single (***kwargs*)

Trims and removes low quality subsections and reads from a single-ended FASTQ file

Parameters

- **fastq_file_in** (*str*) – Location of the input fastq file
- **fastq_file_out** (*str*) – Location of the output fastq file
- **params** (*dict*) – Parameters to use in TrimGalore

Returns Indicator of the success of the function

Return type bool

trimgalore_version (**kwargs)

Trims and removes low quality subsections and reads from a singed-ended FASTQ file

Parameters

- **fastq_file_in** (*str*) – Location of the input fastq file
- **fastq_file_out** (*str*) – Location of the output fastq file
- **params** (*dict*) – Parameters to use in TrimGalore

Returns Indicator of the success of the function

Return type bool

4.2 Indexers

4.2.1 Bowtie 2

class `tool.bowtie_indexer.bowtieIndexerTool` (*configuration=None*)

Tool for running indexers over a genome FASTA file

bowtie2_indexer (**kwargs)

Bowtie2 Indexer

Parameters

- **file_loc** (*str*) – Location of the genome assembly FASTA file
- **idx_loc** (*str*) – Location of the output index file

run (*input_files*, *input_metadata*, *output_files*)

Tool for generating assembly aligner index files for use with the Bowtie 2 aligner

Parameters

- **input_files** (*list*) – List with a single str element with the location of the genome assembly FASTA file
- **metadata** (*list*) –

Returns **array** – First element is a list of the index files. Second element is a list of the matching metadata

Return type list

4.2.2 BSgenome Index

class `tool.forge_bsgenome.bsgenomeTool` (*configuration=None*)

Tool for peak calling for iDamID-seq data

bsgenome_creator (**kwargs)

Make BSgenome index files. Uses an R script that wraps the required code.

Parameters

- **genome** (*str*) –
- **circo_chrom** (*str*) – Comma separated list of chromosome ids that are circular in the genome

- **seed_file_param** (*dict*) – Parameters required for the function to build the seed file
- **genome_2bit** (*str*) –
- **chrom_size** (*str*) –
- **seed_file** (*str*) –
- **bsgenome** (*str*) –

static genome_to_2bit (*genome, genome_2bit*)

Generate the 2bit genome file from a FASTA file

Parameters

- **genome** (*str*) – Location of the FASRA genome file
- **genome_2bit** (*str*) – Location of the 2bit genome file

Returns True if successful, False if not.

Return type bool

static get_chrom_size (*genome_2bit, chrom_size, circ_chrom*)

Generate the chrom.size file and identify the available chromosomes in the 2Bit file.

Parameters

- **genome_2bit** (*str*) – Location of the 2bit genome file
- **chrom_size** (*str*) – Location to save the chrom.size file to
- **circ_chrom** (*list*) – List of chromosomes that are known to be circular

Returns

- *If successful 2 lists – [0] : List of the linear chromosomes in the 2bit file [1] : List of circular chromosomes in the 2bit file*
- *Returns (False, False) if there is an IOError*

run (*input_files, input_metadata, output_files*)

The main function to run iNPS for peak calling over a given BAM file and matching background BAM file.

Parameters

- **input_files** (*list*) – List of input bam file locations where 0 is the bam data file and 1 is the matching background bam file
- **metadata** (*dict*) –

Returns

- **output_files** (*list*) – List of locations for the output files.
- **output_metadata** (*list*) – List of matching metadata dict objects

4.2.3 BS-Seeker2 Indexer

class tool.bs_seeker_indexer.**bssIndexerTool** (*configuration=None*)

Script from BS-Seeker2 for building the index for alignment. In this case it uses Bowtie2.

bss_build_index (***kwargs*)

Function to submit the FASTA file for the reference sequence and build the required index file used by the aligner.

Parameters

- **fasta_file** (*str*) – Location of the genome FASTA file
- **aligner** (*str*) – Aligner to use by BS-Seeker2. Currently only bowtie2 is available in this build
- **aligner_path** (*str*) – Location of the aligners binary file
- **bss_path** – Location of the BS-Seeker2 libraries
- **idx_out** (*str*) – Location of the output compressed index file

Returns **bam_out** – Location of the output bam alignment file

Return type *str*

static **get_bss_index_params** (*params*)

Function to handle to extraction of commandline parameters and formatting them for use in the aligner for BWA ALN

Parameters **params** (*dict*) –

Returns

Return type *list*

run (*input_files*, *input_metadata*, *output_files*)

Tool for indexing the genome assembly using BS-Seeker2. In this case it is using Bowtie2

Parameters

- **input_files** (*list*) – FASTQ file
- **metadata** (*list*) –

Returns **array** – Location of the filtered FASTQ file

Return type *list*

4.2.4 BWA

class `tool.bwa_indexer.bwaIndexerTool` (*configuration=None*)

Tool for running indexers over a genome FASTA file

bwa_indexer (***kwargs*)

BWA Indexer

Parameters

- **file_loc** (*str*) – Location of the genome assembly FASTA file
- **idx_out** (*str*) – Location of the output index file

Returns

Return type *bool*

run (*input_files*, *input_metadata*, *output_files*)

Function to run the BWA over a genome assembly FASTA file to generate the matching index for use with the aligner

Parameters

- **input_files** (*dict*) – List containing the location of the genome assembly FASTA file

- **meta_data** (*dict*) –
- **output_files** (*dict*) – List of output files generated

Returns

- **output_files** (*dict*) –
 - index** [*str*] Location of the index file defined in the input parameters
- **output_metadata** (*dict*) –
 - index** [*Metadata*] Metadata relating to the index file

4.2.5 GEM

class `tool.gem_indexer.gemIndexerTool` (*configuration=None*)
Tool for running indexers over a genome FASTA file

gem_indexer (***kwargs*)
GEM Indexer

Parameters

- **genome_file** (*str*) – Location of the genome assembly FASTA file
- **idx_loc** (*str*) – Location of the output index file

run (*input_files, input_metadata, output_files*)
Tool for generating assembly aligner index files for use with the GEM indexer

Parameters

- **input_files** (*list*) – List with a single *str* element with the location of the genome assembly FASTA file
- **input_metadata** (*list*) –

Returns *array* – First element is a list of the index files. Second element is a list of the matching metadata

Return type *list*

4.2.6 Kallisto

class `tool.kallisto_indexer.kallistoIndexerTool` (*configuration=None*)
Tool for running indexers over a genome FASTA file

kallisto_indexer (***kwargs*)
Kallisto Indexer

Parameters

- **file_loc** (*str*) – Location of the cDNA FASTA file for a genome
- **idx_loc** (*str*) – Location of the output index file

run (*input_files, input_metadata, output_files*)
Tool for generating assembly aligner index files for use with Kallisto

Parameters

- **input_files** (*list*) – FASTA file location will all the cDNA sequences for a given genome

- **input_metadata** (*list*) –

Returns array – First element is a list of the index files. Second element is a list of the matching metadata

Return type list

4.3 Aligners

4.3.1 Bowtie2

class `tool.bowtie_aligner.bowtie2AlignerTool` (*configuration=None*)

Tool for aligning sequence reads to a genome using BWA

bowtie2_aligner_paired (***kwargs*)

Bowtie2 Aligner - Paired End

Parameters

- **genome_file_loc** (*str*) – Location of the genomic fasta
- **read_file_loc1** (*str*) – Location of the FASTQ file
- **read_file_loc2** (*str*) – Location of the FASTQ file
- **bam_loc** (*str*) – Location of the output aligned bam file
- **bt2_1_file** (*str*) – Location of the <genome>.1.bt2 index file
- **bt2_2_file** (*str*) – Location of the <genome>.2.bt2 index file
- **bt2_3_file** (*str*) – Location of the <genome>.3.bt2 index file
- **bt2_4_file** (*str*) – Location of the <genome>.4.bt2 index file
- **bt2_rev1_file** (*str*) – Location of the <genome>.rev.1.bt2 index file
- **bt2_rev2_file** (*str*) – Location of the <genome>.rev.2.bt2 index file
- **aln_params** (*dict*) – Alignment parameters

Returns bam_loc – Location of the output file

Return type str

bowtie2_aligner_single (***kwargs*)

Bowtie2 Aligner - Single End

Parameters

- **genome_file_loc** (*str*) – Location of the genomic fasta
- **read_file_loc1** (*str*) – Location of the FASTQ file
- **bam_loc** (*str*) – Location of the output aligned bam file
- **bt2_1_file** (*str*) – Location of the <genome>.1.bt2 index file
- **bt2_2_file** (*str*) – Location of the <genome>.2.bt2 index file
- **bt2_3_file** (*str*) – Location of the <genome>.3.bt2 index file
- **bt2_4_file** (*str*) – Location of the <genome>.4.bt2 index file
- **bt2_rev1_file** (*str*) – Location of the <genome>.rev.1.bt2 index file

- **bt2_rev2_file** (*str*) – Location of the <genome>.rev.2.bt2 index file
- **aln_params** (*dict*) – Alignment parameters

Returns **bam_loc** – Location of the output file

Return type *str*

static get_aln_params (*params*, *paired=False*)

Function to handle to extraction of commandline parameters and formatting them for use in the aligner for Bowtie2

Parameters

- **params** (*dict*) –
- **paired** (*bool*) – Indicate if the parameters are paired-end specific. [DEFAULT=False]

Returns

Return type *list*

run (*input_files*, *input_metadata*, *output_files*)

The main function to align bam files to a genome using Bowtie2

Parameters

- **input_files** (*dict*) – File 0 is the genome file location, file 1 is the FASTQ file
- **metadata** (*dict*) –
- **output_files** (*dict*) –

Returns

- **output_files** (*dict*) – First element is a list of output_bam_files, second element is the matching meta data
- **output_metadata** (*dict*)

untar_index (***kwargs*)

Extracts the Bowtie2 index files from the genome index tar file.

Parameters

- **genome_file_name** (*str*) – Location string of the genome fasta file
- **genome_idx** (*str*) – Location of the Bowtie2 index file
- **bt2_1_file** (*str*) – Location of the <genome>.1.bt2 index file
- **bt2_2_file** (*str*) – Location of the <genome>.2.bt2 index file
- **bt2_3_file** (*str*) – Location of the <genome>.3.bt2 index file
- **bt2_4_file** (*str*) – Location of the <genome>.4.bt2 index file
- **bt2_rev1_file** (*str*) – Location of the <genome>.rev.1.bt2 index file
- **bt2_rev2_file** (*str*) – Location of the <genome>.rev.2.bt2 index file

Returns Boolean indicating if the task was successful

Return type *bool*

4.3.2 BWA - ALN

class `tool.bwa_aligner.bwaAlignerTool` (*configuration=None*)
 Tool for aligning sequence reads to a genome using BWA

bwa_aligner_paired (***kwargs*)
 BWA ALN Aligner - Paired End

Parameters

- **genome_file_loc** (*str*) – Location of the genomic fasta
- **read_file_loc1** (*str*) – Location of the FASTQ file
- **read_file_loc2** (*str*) – Location of the FASTQ file
- **bam_loc** (*str*) – Location of the output aligned bam file
- **amb_file** (*str*) – Location of the amb index file
- **ann_file** (*str*) – Location of the ann index file
- **bwt_file** (*str*) – Location of the bwt index file
- **pac_file** (*str*) – Location of the pac index file
- **sa_file** (*str*) – Location of the sa index file
- **aln_params** (*dict*) – Alignment parameters

Returns **bam_loc** – Location of the output file

Return type `str`

bwa_aligner_single (***kwargs*)
 BWA ALN Aligner - Single Ended

Parameters

- **genome_file_loc** (*str*) – Location of the genomic fasta
- **read_file_loc** (*str*) – Location of the FASTQ file
- **bam_loc** (*str*) – Location of the output aligned bam file
- **amb_file** (*str*) – Location of the amb index file
- **ann_file** (*str*) – Location of the ann index file
- **bwt_file** (*str*) – Location of the bwt index file
- **pac_file** (*str*) – Location of the pac index file
- **sa_file** (*str*) – Location of the sa index file
- **aln_params** (*dict*) – Alignment parameters

Returns **bam_loc** – Location of the output file

Return type `str`

static **get_aln_params** (*params*)
 Function to handle to extraction of commandline parameters and formatting them for use in the aligner for BWA ALN

Parameters **params** (*dict*) –

Returns

Return type list

run (*input_files*, *input_metadata*, *output_files*)

The main function to align bam files to a genome using BWA

Parameters

- **input_files** (*dict*) – File 0 is the genome file location, file 1 is the FASTQ file
- **metadata** (*dict*) –
- **output_files** (*dict*) –

Returns

- **output_files** (*dict*) – First element is a list of output_bam_files, second element is the matching meta data
- **output_metadata** (*dict*)

untar_index (***kwargs*)

Extracts the BWA index files from the genome index tar file.

Parameters

- **genome_file_name** (*str*) – Location string of the genome fasta file
- **genome_idx** (*str*) – Location of the BWA index file
- **amb_file** (*str*) – Location of the amb index file
- **ann_file** (*str*) – Location of the ann index file
- **bwt_file** (*str*) – Location of the bwt index file
- **pac_file** (*str*) – Location of the pac index file
- **sa_file** (*str*) – Location of the sa index file

Returns Boolean indicating if the task was successful

Return type bool

4.3.3 BWA - MEM

class `tool.bwa_mem_aligner.bwaAlignerMEMTool` (*configuration=None*)

Tool for aligning sequence reads to a genome using BWA

bwa_aligner_paired (***kwargs*)

BWA MEM Aligner - Paired End

Parameters

- **genome_file_loc** (*str*) – Location of the genomic fasta
- **read_file_loc1** (*str*) – Location of the FASTQ file
- **read_file_loc2** (*str*) – Location of the FASTQ file
- **bam_loc** (*str*) – Location of the output aligned bam file
- **amb_file** (*str*) – Location of the amb index file
- **ann_file** (*str*) – Location of the ann index file
- **bwt_file** (*str*) – Location of the bwt index file
- **pac_file** (*str*) – Location of the pac index file

- **sa_file** (*str*) – Location of the sa index file
- **mem_params** (*dict*) – Alignment parameters

Returns **bam_loc** – Location of the output file

Return type *str*

bwa_aligner_single (***kwargs*)

BWA MEM Aligner - Single Ended

Parameters

- **genome_file_loc** (*str*) – Location of the genomic fasta
- **read_file_loc** (*str*) – Location of the FASTQ file
- **bam_loc** (*str*) – Location of the output aligned bam file
- **amb_file** (*str*) – Location of the amb index file
- **ann_file** (*str*) – Location of the ann index file
- **bwt_file** (*str*) – Location of the bwt index file
- **pac_file** (*str*) – Location of the pac index file
- **sa_file** (*str*) – Location of the sa index file
- **mem_params** (*dict*) – Alignment parameters

Returns **bam_loc** – Location of the output file

Return type *str*

static get_mem_params (*params*)

Function to handle to extraction of commandline parameters and formatting them for use in the aligner for BWA MEM

Parameters **params** (*dict*) –

Returns

Return type *list*

run (*input_files, input_metadata, output_files*)

The main function to align bam files to a genome using BWA

Parameters

- **input_files** (*dict*) – File 0 is the genome file location, file 1 is the FASTQ file
- **metadata** (*dict*) –
- **output_files** (*dict*) –

Returns

- **output_files** (*dict*) – First element is a list of output_bam_files, second element is the matching meta data
- **output_metadata** (*dict*)

untar_index (***kwargs*)

Extracts the BWA index files from the genome index tar file.

Parameters

- **genome_file_name** (*str*) – Location string of the genome fasta file

- **genome_idx** (*str*) – Location of the BWA index file
- **amb_file** (*str*) – Location of the amb index file
- **ann_file** (*str*) – Location of the ann index file
- **bwt_file** (*str*) – Location of the bwt index file
- **pac_file** (*str*) – Location of the pac index file
- **sa_file** (*str*) – Location of the sa index file

Returns Boolean indicating if the task was successful

Return type bool

4.3.4 BS-Seeker2 Aligner

class `tool.bs_seeker_aligner.bssAlignerTool` (*configuration=None*)

Script from BS-Seeker2 for building the index for alignment. In this case it uses Bowtie2.

bs_seeker_aligner (***kwargs*)

Alignment of the paired ends to the reference genome

Generates bam files for the alignments

This is performed by running the external program rather than reimplementing the code from the main function to make it easier when it comes to updating the changes in BS-Seeker2

Parameters

- **input_fastq1** (*str*) – Location of paired end FASTQ file 1
- **input_fastq2** (*str*) – Location of paired end FASTQ file 2
- **aligner** (*str*) – Aligner to use
- **aligner_path** (*str*) – Location of the aligner
- **genome_fasta** (*str*) – Location of the genome FASTA file
- **genome_idx** (*str*) – Location of the tar.gz genome index file
- **bam_out** (*str*) – Location of the aligned bam file

Returns **bam_out** – Location of the BAM file generated during the alignment.

Return type file

bs_seeker_aligner_single (***kwargs*)

Alignment of the paired ends to the reference genome

Generates bam files for the alignments

This is performed by running the external program rather than reimplementing the code from the main function to make it easier when it comes to updating the changes in BS-Seeker2

Parameters

- **input_fastq1** (*str*) – Location of paired end FASTQ file 1
- **input_fastq2** (*str*) – Location of paired end FASTQ file 2
- **aligner** (*str*) – Aligner to use
- **aligner_path** (*str*) – Location of the aligner
- **genome_fasta** (*str*) – Location of the genome FASTA file

- **genome_idx** (*str*) – Location of the tar.gz genome index file
- **bam_out** (*str*) – Location of the aligned bam file

Returns **bam_out** – Location of the BAM file generated during the alignment.

Return type file

static get_aln_params (*params*, *paired=False*)

Function to handle to extraction of commandline parameters and formatting them for use in the aligner for Bowtie2

Parameters

- **params** (*dict*) –
- **paired** (*bool*) – Indicate if the parameters are paired-end specific. [DEFAULT=False]

Returns

Return type list

run (*input_files*, *input_metadata*, *output_files*)

Tool for indexing the genome assembly using BS-Seeker2. In this case it is using Bowtie2

Parameters

- **input_files** (*list*) – FASTQ file
- **output_files** (*list*) – Results files.
- **metadata** (*list*) –

Returns **array** – Location of the filtered FASTQ file

Return type list

run_aligner (*genome_idx*, *bam_out*, *script*, *params*)

Run the aligner

Parameters

- **genome_idx** (*str*) – Location of the genome index archive
- **bam_out** (*str*) – Location of the output bam file
- **script** (*str*) – Location of the BS Seeker2 aligner script
- **params** (*list*) – Parameter list for the aligner

Returns True if the function completed successfully

Return type bool

4.4 Filters

4.4.1 BioBamBam Filter

class `tool.biobambam_filter.biobambam` (*configuration=None*)

Tool to sort and filter bam files

biobambam_filter_alignments (***kwargs*)

Sorts and filters the bam file.

It is important that all duplicate alignments have been removed. This can be run as an intermediate step, but should always be run as a check to ensure that the files are sorted and duplicates have been removed.

Parameters

- **bam_file_in** (*str*) – Location of the input bam file
- **bam_file_out** (*str*) – Location of the output bam file
- **tmp_dir** (*str*) – Tmp location for intermediate files during the sorting

Returns **bam_file_out** – Location of the output bam file

Return type *str*

run (*input_files*, *input_metadata*, *output_files*)

The main function to run BioBAMBAMfilter to remove duplicates and spurious reads from the FASTQ files before analysis.

Parameters

- **input_files** (*dict*) – List of input bam file locations where 0 is the bam data file
- **metadata** (*dict*) – Matching meta data for the input files
- **output_files** (*dict*) – List of output file locations

Returns

- **output_files** (*dict*) – Filtered bam file.
- **output_metadata** (*dict*) – List of matching metadata dict objects

4.4.2 BS-Seeker2 Filter

class `tool.bs_seeker_filter.filterReadsTool` (*configuration=None*)

Script from BS-Seeker2 for filtering FASTQ files to remove repeats

bss_seeker_filter (***kwargs*)

This is optional, but removes reads that can be problematic for the alignment of whole genome datasets.

If performing RRBS then this step can be skipped

This is a function that is installed as part of the BS-Seeker installation process.

Parameters **infile** (*str*) – Location of the FASTQ file

Returns **outfile** – Location of the filtered FASTQ file

Return type *str*

run (*input_files*, *input_metadata*, *output_files*)

Tool for filtering duplicate entries from FASTQ files using BS-Seeker2

Parameters

- **input_files** (*list*) – FASTQ file
- **input_metadata** (*list*) –

Returns **array** – Location of the filtered FASTQ file

Return type *list*

4.4.3 Trim Galore

class `tool.trimgalore.trimgalore` (*configuration=None*)

Tool for trimming FASTQ reads that are of low quality

static `get_trimgalore_params` (*params*)

Function to handle for extraction of commandline parameters

Parameters `params` (*dict*) –

Returns

Return type list

run (*input_files, input_metadata, output_files*)

The main function to run TrimGalore to remove low quality and very short reads. TrimGalore uses CutAdapt and FASTQC for the analysis.

Parameters

- **input_files** (*dict*) –

- **fastq1** [string] Location of the FASTQ file

- **fastq2** [string] [OPTIONAL] Location of the paired end FASTQ file

- **metadata** (*dict*) – Matching metadata for the input FASTQ files

Returns

- **output_files** (*dict*) –

- **fastq1_trimmed** [str] Location of the trimmed FASTQ file

- **fastq2_trimmed** [str] [OPTIONAL] Location of a trimmed paired end FASTQ file

- **output_metadata** (*dict*) – Matching metadata for the output files

trimgalore_paired (***kwargs*)

Trims and removes low quality subsections and reads from paired-end FASTQ files

Parameters

- **fastq_file_in** (*str*) – Location of the input fastq file

- **fastq_file_out** (*str*) – Location of the output fastq file

- **params** (*dict*) – Parameters to use in TrimGalore

Returns Indicator of the success of the function

Return type bool

trimgalore_single (***kwargs*)

Trims and removes low quality subsections and reads from a single-ended FASTQ file

Parameters

- **fastq_file_in** (*str*) – Location of the input fastq file

- **fastq_file_out** (*str*) – Location of the output fastq file

- **params** (*dict*) – Parameters to use in TrimGalore

Returns Indicator of the success of the function

Return type bool

`trimgalore_version` (**kwargs)

Trims and removes low quality subsections and reads from a singed-ended FASTQ file

Parameters

- `fastq_file_in` (*str*) – Location of the input fastq file
- `fastq_file_out` (*str*) – Location of the output fastq file
- `params` (*dict*) – Parameters to use in TrimGalore

Returns Indicator of the success of the function

Return type bool

4.5 Peak Calling

4.5.1 BS-Seeker2 Methylation Caller

4.5.2 iDEAR

`class` `tool.idear.idearTool` (*configuration=None*)

Tool for peak calling for iDamID-seq data

`idear_peak_calling` (**kwargs)

Make iDamID-seq peak calls. These are saved as bed files That can then get displayed on genome browsers. Uses an R script that wraps teh iDEAR protocol.

Parameters

- `sample_name` (*str*) –
- `bg_name` (*str*) –
- `sample_bam_tar_file` (*str*) – Location of the aligned sequences in bam format
- `bg_bam_tar_file` (*str*) – Location of the aligned background sequences in bam format
- `species` (*str*) – Species name for the alignments
- `assembly` (*str*) – Assembly used for teh aligned sequences
- `peak_bed` (*str*) – Location of the peak bed file

Returns `peak_bed` – Location of the collated bed file

Return type str

`run` (*input_files, input_metadata, output_files*)

The main function to run iNPS for peak calling over a given BAM file and matching background BAM file.

Parameters

- `input_files` (*list*) – List of input bam file locations where 0 is the bam data file and 1 is the matching background bam file
- `metadata` (*dict*) –

Returns

- `output_files` (*list*) – List of locations for the output files.

- **output_metadata** (*list*) – List of matching metadata dict objects

4.5.3 iNPS

class `tool.inps.inps` (*configuration=None*)

Tool for peak calling for MNase-seq data

inps_peak_calling (***kwargs*)

Convert Bam to Bed then make Nucleosome peak calls. These are saved as bed files That can then get displayed on genome browsers.

Parameters

- **bam_file** (*str*) – Location of the aligned sequences in bam format
- **peak_bed** (*str*) – Location of the collated bed file of nucleosome peak calls

Returns **peak_bed** – Location of the collated bed file of nucleosome peak calls

Return type `str`

run (*input_files, input_metadata, output_files*)

The main function to run iNPS for peak calling over a given BAM file and matching background BAM file.

Parameters

- **input_files** (*list*) – List of input bam file locations where 0 is the bam data file and 1 is the matching background bam file
- **metadata** (*dict*) –

Returns

- **output_files** (*list*) – List of locations for the output files.
- **output_metadata** (*list*) – List of matching metadata dict objects

4.5.4 Kallisto Quantification

class `tool.kallisto_quant.kallistoQuantificationTool` (*configuration=None*)

Tool for quantifying RNA-seq alignments to calculate expression levels of genes within a genome.

kallisto_quant_paired (***kwargs*)

Kallisto quantifier for paired end RNA-seq data

Parameters

- **idx_loc** (*str*) – Location of the output index file
- **fastq_file_loc_01** (*str*) – Location of the FASTQ sequence file
- **fastq_file_loc_02** (*str*) – Location of the paired FASTQ sequence file

Returns **wig_file_loc** – Location of the wig file containing the levels of expression

Return type `loc`

kallisto_quant_single (***kwargs*)

Kallisto quantifier for single end RNA-seq data

Parameters

- **idx_loc** (*str*) – Location of the output index file

- **fastq_file_loc** (*str*) – Location of the FASTQ sequence file

Returns **wig_file_loc** – Location of the wig file containing the levels of expression

Return type loc

kallisto_tsv2bed (***kwargs*)

So that the TSV file can be viewed within the genome browser it is handy to convert the file to a BigBed file

kallisto_tsv2gff (***kwargs*)

So that the TSV file can be viewed within the genome browser it is handy to convert the file to a BigBed file

static load_gff_ensembl (*gff_file*)

Function to extract all of the genes and their locations from a GFF file generated by ensembl

static load_gff_ucsc (*gff_file*)

Function to extract all of the genes and their locations from a GFF file generated by ensembl

run (*input_files, input_metadata, output_files*)

Tool for calculating the level of expression

Parameters

- **input_files** (*list*) – Kallisto index file for the FASTQ file for the experimtnal alignments
- **input_metadata** (*list*) –

Returns **array** – First element is a list of the index files. Second element is a list of the matching metadata

Return type list

static seq_read_stats (*file_in*)

Calculate the mean and standard deviation of the reads in a fastq file

Parameters **file_in** (*str*) – Location of a FASTQ file

Returns mean : Mean length of sequenced strands std : Standard deviation of lengths of sequenced strands

Return type dict

4.5.5 MACS2

class `tool.macs2.macs2` (*configuration=None*)

Tool for peak calling for ChIP-seq data

static get_macs2_params (*params*)

Function to handle to extraction of commandline parameters and formatting them for use in the aligner for BWA ALN

Parameters **params** (*dict*) –

Returns

Return type list

macs2_peak_calling (***kwargs*)

Function to run MACS2 for peak calling on aligned sequence files and normalised against a provided background set of alignments.

Parameters

- **name** (*str*) – Name to be used to identify the files
- **bam_file** (*str*) – Location of the aligned FASTQ files as a bam file
- **bai_file** (*str*) – Location of the bam index file
- **bam_file_bgd** (*str*) – Location of the aligned FASTQ files as a bam file representing background values for the cell
- **bai_file_bgd** (*str*) – Location of the background bam index file
- **narrowpeak** (*str*) – Location of the output narrowpeak file
- **summits_bed** (*str*) – Location of the output summits bed file
- **broadpeak** (*str*) – Location of the output broadpeak file
- **gappedpeak** (*str*) – Location of the output gappedpeak file
- **chromosome** (*str*) – If the tool is to be run over a single chromosome the matching chromosome name should be specified. If None then the whole bam file is analysed

Returns

- **narrowPeak** (*file*) – BED6+4 file - ideal for transcription factor binding site identification
- **summitPeak** (*file*) – BED4+1 file - Contains the peak summit locations for every peak
- **broadPeak** (*file*) – BED6+3 file - ideal for histone binding site identification
- **gappedPeak** (*file*) – BED12+3 file - Contains a merged set of the broad and narrow peak files
- *Definitions defined for each of these files have come from the MACS2*
- **documentation described in the docs at [https \(//github.com/taoliu/MACS\)](https://github.com/taoliu/MACS)**

macs2_peak_calling_nobgd (***kwargs*)

Function to run MACS2 for peak calling on aligned sequence files without a background dataset for normalisation.

Parameters

- **name** (*str*) – Name to be used to identify the files
- **bam_file** (*str*) – Location of the aligned FASTQ files as a bam file
- **bai_file** (*str*) – Location of the bam index file
- **narrowpeak** (*str*) – Location of the output narrowpeak file
- **summits_bed** (*str*) – Location of the output summits bed file
- **broadpeak** (*str*) – Location of the output broadpeak file
- **gappedpeak** (*str*) – Location of the output gappedpeak file
- **chromosome** (*str*) – If the tool is to be run over a single chromosome the matching chromosome name should be specified. If None then the whole bam file is analysed

Returns

- **narrowPeak** (*file*) – BED6+4 file - ideal for transcription factor binding site identification
- **summitPeak** (*file*) – BED4+1 file - Contains the peak summit locations for every peak
- **broadPeak** (*file*) – BED6+3 file - ideal for histone binding site identification

- **gappedPeak** (*file*) – BED12+3 file - Contains a merged set of the broad and narrow peak files
- *Definitions defined for each of these files have come from the MACS2*
- **documentation described in the docs at [https](https://github.com/taoliu/MACS)** (*//github.com/taoliu/MACS*)

run (*input_files, input_metadata, output_files*)

The main function to run MACS 2 for peak calling over a given BAM file and matching background BAM file.

Parameters

- **input_files** (*dict*) – List of input bam file locations where 0 is the bam data file and 1 is the matching background bam file
- **metadata** (*dict*) –

Returns

- **output_files** (*dict*) – List of locations for the output files.
- **output_metadata** (*dict*) – List of matching metadata dict objects

4.6 Hi-C Parsing

The following tools are a split out of the Hi-C pipelines generated to use the TADbit library.

4.6.1 FASTQ mapping

class `tool.tb_full_mapping.tbFullMappingTool`

Tool for mapping fastq paired end files to the GEM index files

run (*input_files, input_metadata, output_files*)

The main function to map the FASTQ files to the GEM file over different window sizes ready for alignment

Parameters

- **input_files** (*list*) –
 - gem_file** [str] Location of the genome GEM index file
 - fastq_file_bgd** [str] Location of the FASTQ file
- **metadata** (*dict*) –
 - windows** [list] List of lists with the window sizes to be computed
 - enzyme_name** [str] Restriction enzyme used [OPTIONAL]

Returns

- **output_files** (*list*) – List of locations for the output files.
- **output_metadata** (*list*) – List of matching metadata dict objects

tb_full_mapping_frag (***kwargs*)

Function to map the FASTQ files to the GEM file based on fragments derived from the restriction enzyme that was used.

Parameters

- **gem_file** (*str*) – Location of the genome GEM index file

- **fastq_file_bgd** (*str*) – Location of the FASTQ file
- **enzyme_name** (*str*) – Restriction enzyme name (MboI)
- **windows** (*list*) – List of lists with the window sizes to be computed
- **window_file** (*str*) – Location of the first window index file

Returns **window_file** – Location of the window index file

Return type *str*

tb_full_mapping_iter (***kwargs*)

Function to map the FASTQ files to the GEM file over different window sizes ready for alignment

Parameters

- **gem_file** (*str*) – Location of the genome GEM index file
- **fastq_file_bgd** (*str*) – Location of the FASTQ file
- **windows** (*list*) – List of lists with the window sizes to be computed
- **window1** (*str*) – Location of the first window index file
- **window2** (*str*) – Location of the second window index file
- **window3** (*str*) – Location of the third window index file
- **window4** (*str*) – Location of the fourth window index file

Returns

- **window1** (*str*) – Location of the first window index file
- **window2** (*str*) – Location of the second window index file
- **window3** (*str*) – Location of the third window index file
- **window4** (*str*) – Location of the fourth window index file

4.6.2 Map Parsing

class `tool.tb_parse_mapping.tbParseMappingTool`

Tool for parsing the mapped reads and generating the list of paired ends that have a match at both ends.

run (*input_files, input_metadata, output_files*)

The main function to map the aligned reads and return the matching pairs. Parsing of the mappings can be either iterative or fragment based. If it is to be iterative then the locations of 4 output file windows for each end of the paired end window need to be provided. If it is fragment based, then only 2 window locations need to be provided along with an enzyme name.

Parameters

- **input_files** (*list*) –
 - genome_file** [str] Location of the genome FASTA file
 - window1_1** [str] Location of the first window index file
 - window1_2** [str] Location of the second window index file
 - window1_3** [str] [OPTIONAL] Location of the third window index file
 - window1_4** [str] [OPTIONAL] Location of the fourth window index file
 - window2_1** [str] Location of the first window index file

window2_2 [str] Location of the second window index file

window2_3 [str] [OPTIONAL] Location of the third window index file

window2_4 [str] [OPTIONAL] Location of the fourth window index file

- **metadata** (*dict*) –

windows [list] List of lists with the window sizes to be computed

enzyme_name [str] Restricture enzyme name

mapping [list] The mapping function used. The options are iter or frag.

Returns

- **output_files** (*list*) – List of locations for the output files.

- **output_metadata** (*dict*) – Dict of matching metadata dict objects

Example

Iterative:

```
from tool import tb_parse_mapping

genome_file = 'genome.fasta'

root_name_1 = "/tmp/data/expt_source_1".split
root_name_2 = "/tmp/data/expt_source_2".split
windows = [[1,25], [1,50], [1,75], [1,100]]

windows1 = []
windows2 = []

for w in windows:
    tail = "_full_" + w[0] + "-" + w[1] + ".map"
    windows1.append('/'.join(root_name_1) + tail)
    windows2.append('/'.join(root_name_2) + tail)

files = [genome_file] + windows1 + windows2

tpm = tb_parse_mapping.tb_parse_mapping()
metadata = {'enzyme_name' : 'MboI', 'mapping' : ['iter', 'iter'], 'expt_name' :
↳ 'test'}
tpm_files, tpm_meta = tpm.run(files, metadata)
```

Fragment based mapping:

```
from tool import tb_parse_mapping

genome_file = 'genome.fasta'

root_name_1 = "/tmp/data/expt_source_1".split
root_name_2 = "/tmp/data/expt_source_2".split
windows = [[1,100]]

start = windows[0][0]
end = windows[0][1]
```

(continues on next page)

(continued from previous page)

```

window1_1 = '/'.join(root_name_1) + "_full_" + start + "-" + end + ".map"
window1_2 = '/'.join(root_name_1) + "_frag_" + start + "-" + end + ".map"

window2_1 = '/'.join(root_name_2) + "_full_" + start + "-" + end + ".map"
window2_2 = '/'.join(root_name_2) + "_frag_" + start + "-" + end + ".map"

files = [
    genome_file,
    window1_1, window1_2,
    window2_1, window2_2,
]

tpm = tb_parse_mapping.tb_parse_mapping()
metadata = {'enzyme_name' : 'MboI', 'mapping' : ['frag', 'frag'], 'expt_name' :
↳='test'}
tpm_files, tpm_meta = tpm.run(files, metadata)

```

tb_parse_mapping_frag (**kwargs)

Function to map the aligned reads and return the matching pairs

Parameters

- **genome_seq** (*dict*) – Object containing the sequence of each of the chromosomes
- **enzyme_name** (*str*) – Name of the enzyme used to digest the genome
- **window1_full** (*str*) – Location of the first window index file
- **window1_frag** (*str*) – Location of the second window index file
- **window2_full** (*str*) – Location of the first window index file
- **window2_frag** (*str*) – Location of the second window index file
- **reads** (*str*) – Location of the reads that has a matching location at both ends of the paired reads

Returns **reads** – Location of the intersection of mapped reads that have matching reads in both pair end files

Return type *str*

tb_parse_mapping_iter (**kwargs)

Function to map the aligned reads and return the matching pairs

Parameters

- **genome_seq** (*dict*) – Object containing the sequence of each of the chromosomes
- **enzyme_name** (*str*) – Name of the enzyme used to digest the genome
- **window1_1** (*str*) – Location of the first window index file
- **window1_2** (*str*) – Location of the second window index file
- **window1_3** (*str*) – Location of the third window index file
- **window1_4** (*str*) – Location of the fourth window index file
- **window2_1** (*str*) – Location of the first window index file
- **window2_2** (*str*) – Location of the second window index file
- **window2_3** (*str*) – Location of the third window index file

- **window2_4** (*str*) – Location of the fourth window index file
- **reads** (*str*) – Location of the reads that has a matching location at both ends of the paired reads

Returns **reads** – Location of the intersection of mapped reads that have matching reads in both pair end files

Return type `str`

4.6.3 Filter Aligned Reads

class `tool.tb_filter.tbFilterTool` (*configuration=None*)

Tool for filtering out experimental artifacts from the aligned data

run (*input_files, input_metadata, output_files*)

The main function to filter the reads to remove experimental artifacts

Parameters

- **input_files** (*list*) –
- **reads** [*str*] Location of the reads that has a matching location at both ends of the paired reads
- **metadata** (*dict*) –
- **conservative** [*bool*] Level of filtering to apply [DEFAULT : True]

Returns

- **output_files** (*list*) – List of locations for the output files.
- **output_metadata** (*list*) – List of matching metadata dict objects

tb_filter (***kwargs*)

Function to filter out experimental artifacts

Parameters

- **reads** (*str*) – Location of the reads that has a matching location at both ends of the paired reads
- **filtered_reads_file** (*str*) – Location of the filtered reads
- **conservative** (*bool*) – Level of filtering [DEFAULT : True]

Returns **filtered_reads** – Location of the filtered reads

Return type `str`

4.6.4 Identify TADs and Compartments

class `tool.tb_segment.tbSegmentTool`

Tool for finding tads and compartments in an adjacency matrix

run (*input_files, input_metadata, output_files*)

The main function to the predict TAD sites and compartments for a given resolution from the Hi-C matrix

Parameters

- **input_files** (*list*) –
- **bamin** [*str*] Location of the tadbit bam paired reads

biases [str] Location of the pickle hic biases

- **metadata** (*dict*) –

resolution [int] Resolution of the Hi-C

workdir [str] Location of working directory

ncpus [int] Number of cpus to use

Returns

- **output_files** (*list*) – List of locations for the output files.
- **output_metadata** (*list*) – List of matching metadata dict objects

tb_segment (***kwargs*)

Function to find tads and compartments in the Hi-C matrix

Parameters

- **bamin** (*str*) – Location of the tadbit bam paired reads
- **biases** (*str*) – Location of the pickle hic biases
- **resolution** (*int*) – Resolution of the Hi-C
- **callers** (*str*) – 1 for ta calling, 2 for compartment calling
- **workdir** (*str*) – Location of working directory
- **ncpus** (*int*) – Number of cpus to use

Returns

- **compartments** (*str*) – Location of tsv file with compartment definition
- **tads** (*str*) – Location of tsv file with tad definition
- **filtered_bins** (*str*) – Location of filtered_bins png

4.6.5 Normalize paired end reads file

class `tool.tb_normalize.tbNormalizeTool`

Tool for normalizing an adjacency matrix

run (*input_files, input_metadata, output_files*)

The main function for the normalization of the Hi-C matrix to a given resolution

Parameters

- **input_files** (*list*) –
- **bamin** [str] Location of the tadbit bam paired reads
- **metadata** (*dict*) –
- **normalization: str** normalization(s) to apply. Order matters. Choices: [Vanilla, oneD]
- **resolution** [str] Resolution of the Hi-C
- **min_perc** [str] lower percentile from which consider bins as good.
- **max_perc** [str] upper percentile until which consider bins as good.
- **workdir** [str] Location of working directory
- **ncpus** [str] Number of cpus to use

min_count [str] minimum number of reads mapped to a bin (recommended value could be 2500). If set this option overrides the `perc_zero`

fasta: str Location of the fasta file with genome sequence, to compute GC content and number of restriction sites per bin. Required for oneD normalization

mappability: str Location of the file with mappability, required for oneD normalization

rest_enzyme: str For oneD normalization. Name of the restriction enzyme used to do the Hi-C experiment

Returns

- **output_files** (*list*) – List of locations for the output files.
- **output_metadata** (*list*) – List of matching metadata dict objects

tb_normalize (**kwargs)

Function to normalize to a given resolution the Hi-C matrix

Parameters

- **bamin** (*str*) – Location of the tadbit bam paired reads
- **normalization** (*str*) – normalization(s) to apply. Order matters. Choices: [Vanilla, oneD]
- **resolution** (*str*) – Resolution of the Hi-C
- **min_perc** (*str*) – lower percentile from which consider bins as good.
- **max_perc** (*str*) – upper percentile until which consider bins as good.
- **workdir** (*str*) – Location of working directory
- **ncpus** (*str*) – Number of cpus to use
- **min_count** (*str*) – minimum number of reads mapped to a bin (recommended value could be 2500). If set this option overrides the `perc_zero`
- **fasta** (*str*) – Location of the fasta file with genome sequence, to compute GC content and number of restriction sites per bin. Required for oneD normalization
- **mappability** (*str*) – Location of the file with mappability, required for oneD normalization
- **rest_enzyme** (*str*) – For oneD normalization. Name of the restriction enzyme used to do the Hi-C experiment

Returns

- **hic_biases** (*str*) – Location of HiC biases pickle file
- **interactions** (*str*) – Location of interaction decay vs genomic distance pdf
- **filtered_bins** (*str*) – Location of filtered_bins png

4.6.6 Extract binned matrix from paired end reads file

```
class tool.tb_bin.tbBinTool
```

```
    Tool for binning an adjacency matrix
```

```
    run (input_files, input_metadata, output_files)
```

```
        The main function to the predict TAD sites for a given resolution from the Hi-C matrix
```

Parameters

- **input_files** (*list*) –
 - bamin** [str] Location of the tadbit bam paired reads
 - biases** [str] Location of the pickle hic biases
- **input_metadata** (*dict*) –
 - resolution** [int] Resolution of the Hi-C
 - coord1** [str] Coordinate of the region to retrieve. By default all genome, arguments can be either one chromosome name, or the coordinate in the form: “-c chr3:110000000-120000000”
 - coord2** [str] Coordinate of a second region to retrieve the matrix in the intersection with the first region.
 - norm** [str] [['raw']] normalization(s) to apply. Order matters. Choices: [norm, decay, raw]
 - workdir** [str] Location of working directory
 - ncpus** [int] Number of cpus to use

Returns

- **output_files** (*list*) – List of locations for the output files.
- **output_metadata** (*list*) – List of matching metadata dict objects

tb_bin (***kwargs*)

Function to bin to a given resolution the Hi-C matrix

Parameters

- **bamin** (*str*) – Location of the tadbit bam paired reads
- **biases** (*str*) – Location of the pickle hic biases
- **resolution** (*int*) – Resolution of the Hi-C
- **coord1** (*str*) – Coordinate of the region to retrieve. By default all genome, arguments can be either one chromosome name, or the coordinate in the form: “-c chr3:110000000-120000000”
- **coord2** (*str*) – Coordinate of a second region to retrieve the matrix in the intersection with the first region.
- **norm** (*list*) – [['raw']] normalization(s) to apply. Order matters. Choices: [norm, decay, raw]
- **workdir** (*str*) – Location of working directory
- **ncpus** (*int*) – Number of cpus to use

Returns

- **hic_contacts_matrix_raw** (*str*) – Location of HiC raw matrix in text format
- **hic_contacts_matrix_nrm** (*str*) – Location of HiC normalized matrix in text format
- **hic_contacts_matrix_raw_fig** (*str*) – Location of HiC raw matrix in png format
- **hic_contacts_matrix_norm_fig** (*str*) – Location of HiC normalized matrix in png format

4.6.7 Save Matrix to HDF5 File

class `tool.tb_save_hdf5_matrix.tbSaveAdjacencyHDF5Tool`

Tool for filtering out experimental artifacts from the aligned data

run (*input_files*, *output_files*, *metadata=None*)

The main function save the adjacency list from Hi-C into an HDF5 index file at the defined resolutions.

Parameters

- **input_files** (*list*) –
adj_list [str] Location of the adjacency list
hdf5_file [str] Location of the HDF5 output matrix file
- **metadata** (*dict*) –
resolutions [list] Levels of resolution for the adjacency list to be saved at
assembly [str] Assembly of the aligned sequences
normalized [bool] Whether the dataset should be normalised before saving

Returns

- **output_files** (*list*) – List of locations for the output files.
- **output_metadata** (*list*) – List of matching metadata dict objects

tb_matrix_hdf5 (***kwargs*)

Function to the Hi-C matrix into an HDF5 file

This has to be run sequentially as it is not possible for multiple streams to write to the same HDF5 file. This is a run once and leave operation. There also needs to be a check that no other process is writing to the HDF5 file at the same time. This should be done at the stage and unstaging level to prevent to file getting written to by multiple processes and generating conflicts.

This needs to include attributes for the chromosomes for each resolution - See the `mg-rest-adjacency-hdf5-reader` for further details about the requirement. This prevents the need for secondary storage details outside of the HDF5 file.

Parameters

- **hic_data** (*hic_data*) – Hi-C data object
- **hdf5_file** (*str*) – Location of the HDF5 output matrix file
- **resolution** (*int*) – Resolution to read the Hi-C adjacency list at
- **chromosomes** (*list*) – List of lists of the chromosome names and their size in the order that they are presented for indexing

Returns **hdf5_file** – Location of the HDF5 output matrix file

Return type str

4.6.8 Generate TAD Predictions

class `tool.tb_generate_tads.tbGenerateTADsTool`

Tool for taking the adjacency lists and predicting TADs

run (*input_files*, *output_files*, *metadata=None*)

The main function to the predict TAD sites for a given resolution from the Hi-C matrix

Parameters

- **input_files** (*list*) –
adj_list [str] Location of the adjacency list
- **metadata** (*dict*) –
resolutions [list] Levels of resolution for the adjacency list to be daved at
assembly [str] Assembly of the aligned sequences

Returns

- **output_files** (*list*) – List of locations for the output files.
- **output_metadata** (*list*) – List of matching metadata dict objects

tb_generate_tads (***kwargs*)

Function to the predict TAD sites for a given resolution from the Hi-C matrix

Parameters

- **expt_name** (*str*) – Location of the adjacency list
- **matrix_file** (*str*) – Location of the HDF5 output matrix file
- **resolution** (*int*) – Resolution to read the Hi-C adjacency list at
- **tad_file** (*str*) – Location of the output TAD file

Returns **tad_file** – Location of the output TAD file

Return type str

tb_hic_chr (***kwargs*)

Get the list of chromosomes in the adjacency list

tb_merge_tad_files (***kwargs*)

Merge 2 TAD adjacnechy list files

4.6.9 Generate 3D models from binned interaction matrix

class `tool.tb_model.tbModelTool`

Tool for normalizing an adjacency matrix

run (*input_files, input_metadata, output_files*)

The main function for the normalization of the Hi-C matrix to a given resolution

Parameters

- **input_files** (*list*) –
hic_contacts_matrix_norm [str] Location of the tab-separated normalized matrix
- **metadata** (*dict*) –
optimize_only: bool True if only optimize, False for computing the models and stats
gen_pos_chrom_name [str] Coordinates of the genomic region to model.
resolution [str] Resolution of the Hi-C
gen_pos_begin [int] Genomic coordinate from which to start modeling.
gen_pos_end [int] Genomic coordinate where to end modeling.
num_mod_comp [int] Number of models to compute for each optimization step.

- num_mod_comp** [int] Number of models to keep.
- max_dist** [str] Range of numbers for optimal maxdist parameter, i.e. 400:1000:100; or just a single number e.g. 800; or a list of numbers e.g. 400 600 800 1000.
- upper_bound** [int] Range of numbers for optimal upfreq parameter, i.e. 0:1.2:0.3; or just a single number e.g. 0.8; or a list of numbers e.g. 0.1 0.3 0.5 0.9.
- lower_bound** [int] Range of numbers for optimal low parameter, i.e. -1.2:0:0.3; or just a single number e.g. -0.8; or a list of numbers e.g. -0.1 -0.3 -0.5 -0.9.
- cutoff** [str] Range of numbers for optimal cutoff distance. Cutoff is computed based on the resolution. This cutoff distance is calculated taking as reference the diameter of a modeled particle in the 3D model. i.e. 1.5:2.5:0.5; or just a single number e.g. 2; or a list of numbers e.g. 2 2.5.
- workdir** [str] Location of working directory
- ncpus** [str] Number of cpus to use

Returns

- **output_files** (*list*) – List of locations for the output files.
- **output_metadata** (*list*) – List of matching metadata dict objects

tb_model (**kwargs)

Function to normalize to a given resolution the Hi-C matrix

Parameters

- **optimize_only** (*bool*) – True if only optimize, False for computing the models and stats
- **hic_contacts_matrix_norm** (*str*) – Location of the tab-separated normalized matrix
- **resolution** (*str*) – Resolution of the Hi-C
- **gen_pos_chrom_name** (*str*) – Coordinates of the genomic region to model.
- **gen_pos_begin** (*int*) – Genomic coordinate from which to start modeling.
- **gen_pos_end** (*int*) – Genomic coordinate where to end modeling.
- **num_mod_comp** (*int*) – Number of models to compute for each optimization step.
- **num_mod_comp** – Number of models to keep.
- **max_dist** (*str*) – Range of numbers for optimal maxdist parameter, i.e. 400:1000:100; or just a single number e.g. 800; or a list of numbers e.g. 400 600 800 1000.
- **upper_bound** (*int*) – Range of numbers for optimal upfreq parameter, i.e. 0:1.2:0.3; or just a single number e.g. 0.8; or a list of numbers e.g. 0.1 0.3 0.5 0.9.
- **lower_bound** (*int*) – Range of numbers for optimal low parameter, i.e. -1.2:0:0.3; or just a single number e.g. -0.8; or a list of numbers e.g. -0.1 -0.3 -0.5 -0.9.
- **cutoff** (*str*) – Range of numbers for optimal cutoff distance. Cutoff is computed based on the resolution. This cutoff distance is calculated taking as reference the diameter of a modeled particle in the 3D model. i.e. 1.5:2.5:0.5; or just a single number e.g. 2; or a list of numbers e.g. 2 2.5.
- **workdir** (*str*) – Location of working directory
- **ncpus** (*str*) – Number of cpus to use

Returns

- **tadkit_models** (*str*) – Location of TADkit json file
- **modeling_stats** (*str*) – Location of the folder with the modeling files and stats

5.1 Common Functions

The following functions are ones that have been used across multiple tools for transformations of the data when required.

```
class tool.common.cd(newpath)  
    Context manager for changing the current working directory
```

5.2 Alignment Utilities

```
class tool.aligner_utils.alignerUtils
```

Functions for downloading and processing N-seq FastQ files. Functions provided allow for the downloading and indexing of the genome assemblies.

```
static bowtie2_align_reads(genome_file, bam_loc, params, reads_file_1,  
                           reads_file_2=None)
```

Map the reads to the genome using BWA.

Parameters

- **genome_file** (*str*) – Location of the assembly file in the file system
- **reads_file** (*str*) – Location of the reads file in the file system
- **bam_loc** (*str*) – Location of the output file

```
bowtie2_untar_index(genome_name, tar_file, bt2_1_file, bt2_2_file, bt2_3_file, bt2_4_file,  
                   bt2_rev1_file, bt2_rev2_file)
```

Extracts the BWA index files from the genome index tar file.

Parameters

- **genome_file_name** (*str*) – Location string of the genome fasta file
- **tar_file** (*str*) – Location of the Bowtie2 index file

- **bt2_1_file** (*str*) – Location of the amb index file
- **bt2_2_file** (*str*) – Location of the ann index file
- **bt2_3_file** (*str*) – Location of the bwt index file
- **bt2_4_file** (*str*) – Location of the pac index file
- **bt2_rev1_file** (*str*) – Location of the sa index file
- **bt2_rev2_file** (*str*) – Location of the sa index file

Returns Boolean indicating if the task was successful

Return type bool

static bowtie_index_genome (*genome_file*)

Create an index of the genome FASTA file with Bowtie2. These are saved alongside the assembly file.

Parameters **genome_file** (*str*) – Location of the assembly file in the file system

bwa_aln_align_reads_paired (*genome_file, reads_file_1, reads_file_2, bam_loc, params*)

Map the reads to the genome using BWA.

Parameters

- **genome_file** (*str*) – Location of the assembly file in the file system
- **reads_file** (*str*) – Location of the reads file in the file system
- **bam_loc** (*str*) – Location of the output file

bwa_aln_align_reads_single (*genome_file, reads_file, bam_loc, params*)

Map the reads to the genome using BWA. :param genome_file: Location of the assembly file in the file system :type genome_file: str :param reads_file: Location of the reads file in the file system :type reads_file: str :param bam_loc: Location of the output file :type bam_loc: str

static bwa_index_genome (*genome_file*)

Create an index of the genome FASTA file with BWA. These are saved alongside the assembly file. If the index has already been generated then the locations of the files are returned

Parameters **genome_file** (*str*) – Location of the assembly file in the file system

Returns

- **amb_file** (*str*) – Location of the amb file
- **ann_file** (*str*) – Location of the ann file
- **bwt_file** (*str*) – Location of the bwt file
- **pac_file** (*str*) – Location of the pac file
- **sa_file** (*str*) – Location of the sa file

Example

```

1 from tool.aligner_utils import alignerUtils
2 au_handle = alignerUtils()
3
4 indexes = au_handle.bwa_index_genome('/<data_dir>/human_GRCh38.fa.gz')
5 print(indexes)

```

static bwa_mem_align_reads (*genome_file*, *bam_loc*, *params*, *reads_file_1*,
reads_file_2=None)

Map the reads to the genome using BWA.

Parameters

- **genome_file** (*str*) – Location of the assembly file in the file system
- **reads_file** (*str*) – Location of the reads file in the file system
- **bam_loc** (*str*) – Location of the output file

bwa_untar_index (*genome_name*, *tar_file*, *amb_file*, *ann_file*, *bwt_file*, *pac_file*, *sa_file*)

Extracts the BWA index files from the genome index tar file.

Parameters

- **genome_file_name** (*str*) – Location string of the genome fasta file
- **genome_idx** (*str*) – Location of the BWA index file
- **amb_file** (*str*) – Location of the amb index file
- **ann_file** (*str*) – Location of the ann index file
- **bwt_file** (*str*) – Location of the bwt index file
- **pac_file** (*str*) – Location of the pac index file
- **sa_file** (*str*) – Location of the sa index file

Returns Boolean indicating if the task was successful

Return type bool

static gem_index_genome (*genome_file*, *index_name=None*)

Create an index of the genome FASTA file with GEM. These are saved alongside the assembly file.

Parameters **genome_file** (*str*) – Location of the assembly file in the file system

static replaceENAHeader (*file_path*, *file_out*)

The ENA header has pipes in the header as part of the stable_id. This function removes the ENA stable_id and replaces it with the final section after splitting the stable ID on the pipe.

5.3 Bam Utilities

class `tool.bam_utils.bamUtils`

Tool for handling bam files

static bam_copy (*bam_in*, *bam_out*)

Wrapper function to copy from one bam file to another

Parameters

- **bam_in** (*str*) – Location of the input bam file
- **bam_out** (*str*) – Location of the output bam file

static bam_count_reads (*bam_file*, *aligned=False*)

Wrapper to count the number of (aligned) reads in a bam file

static bam_filter (*bam_file*, *bam_file_out*, *filter_name*)

Wrapper for filtering out reads from a bam file

Parameters

- **bam_file** (*str*) –
- **bam_file_out** (*str*) –
- **filter** (*str*) –

One of: duplicate - Read is PCR or optical duplicate (1024) supplementary - Reads that are chimeric, fusion or non linearly aligned (2048) unmapped - Read is unmapped or not the primary alignment (260)

static bam_index (*bam_file, bam_idx_file*)

Wrapper for the pysam SAMtools index function

Parameters

- **bam_file** (*str*) – Location of the bam file that is to be indexed
- **bam_idx_file** (*str*) – Location of the bam index file (.bai)

static bam_list_chromosomes (*bam_file*)

Wrapper to list the chromosome names that are present within the bam file

Parameters **bam_file** (*str*) – Location of the bam file

Returns List of the names of the chromosomes that are present in the bam file

Return type list

static bam_merge (**args*)

Wrapper for the pysam SAMtools merge function

Parameters

- **bam_file_1** (*str*) – Location of the bam file to merge into
- **bam_file_2** (*str*) – Location of the bam file that is to get merged into bam_file_1

static bam_paired_reads (*bam_file*)

Wrapper to test if a bam file contains paired end reads

static bam_sort (*bam_file*)

Wrapper for the pysam SAMtools sort function

Parameters **bam_file** (*str*) – Location of the bam file to sort

static bam_split (*bam_file_in, bai_file, chromosome, bam_file_out*)

Wrapper to extract a single chromosomes worth of reading into a new bam file

Parameters

- **bam_file_in** (*str*) – Location of the input bam file
- **bai_file** (*str*) – Location of the bam index file. This needs to be in the same directory as the bam_file_in
- **chromosome** (*str*) – Name of the chromosome whose alignments are to be extracted
- **bam_file_out** (*str*) – Location of the output bam file

static bam_stats (*bam_file*)

Wrapper for the pysam SAMtools flagstat function

Parameters **bam_file** (*str*) – Location of the bam file

Returns list – qc_passed : int qc_failed : int description : str

Return type dict

static bam_to_bed (*bam_file*, *bed_file*)
Function for converting bam files to bed files

static check_header (*bam_file*)
Wrapper for the pysam SAMtools for checking if a bam file is sorted

Parameters **bool** – True if the file has been sorted

static sam_to_bam (*sam_file*, *bam_file*)
Function for converting sam files to bam files

class `tool.bam_utils.bamUtilsTask`

Wrappers so that the function above can be used as part of a @task within COMPSs avoiding the files being copied around the infrastructure too many times

bam_copy (**kwargs)
Wrapper function to copy from one bam file to another

Parameters

- **bam_in** (*str*) – Location of the input bam file
- **bam_out** (*str*) – Location of the output bam file

bam_filter (**kwargs)
Wrapper for filtering out reads from a bam file

Parameters

- **bam_file** (*str*) –
- **bam_file_out** (*str*) –
- **filter** (*str*) –

One of: duplicate - Read is PCR or optical duplicate (1024) unmapped - Read is unmapped or not the primary alignment (260)

bam_index (**kwargs)
Wrapper for the pysam SAMtools merge function

Parameters

- **bam_file** (*str*) – Location of the bam file that is to be indexed
- **bam_idx_file** (*str*) – Location of the bam index file (.bai)

bam_list_chromosomes (**kwargs)
Wrapper to get the list of chromosomes in a given bam file

Parameters **bam_file** (*str*) – Location of the bam file

Returns **chromosome_list** – List of the chromosomes in the bam file

Return type list

bam_merge (*in_bam_job_files*)
Wrapper task taking any number of bam files and merging them into a single bam file.

Parameters **bam_job_files** (*list*) – List of the locations of the separate bam files that are to be merged The first file in the list will be taken as the output file name

bam_merge_10 (**kwargs)
Wrapper for the pysam SAMtools merge function

Parameters

- **bam_file_1** (*str*) – Location of the bam file to merge into
- **bam_file_2** (*str*) – Location of the bam file that is to get merged into bam_file_1
- **bam_file_3** (*str*) – Location of the bam file that is to get merged into bam_file_1
- **bam_file_4** (*str*) – Location of the bam file that is to get merged into bam_file_1
- **bam_file_5** (*str*) – Location of the bam file that is to get merged into bam_file_1
- **bam_file_6** (*str*) – Location of the bam file that is to get merged into bam_file_1
- **bam_file_7** (*str*) – Location of the bam file that is to get merged into bam_file_1
- **bam_file_8** (*str*) – Location of the bam file that is to get merged into bam_file_1
- **bam_file_9** (*str*) – Location of the bam file that is to get merged into bam_file_1
- **bam_file_10** (*str*) – Location of the bam file that is to get merged into bam_file_1

bam_merge_2 (***kwargs*)

Wrapper for the pysam SAMtools merge function

Parameters

- **bam_file_1** (*str*) – Location of the bam file to merge into
- **bam_file_2** (*str*) – Location of the bam file that is to get merged into bam_file_1

bam_merge_3 (***kwargs*)

Wrapper for the pysam SAMtools merge function

Parameters

- **bam_file_1** (*str*) – Location of the bam file to merge into
- **bam_file_2** (*str*) – Location of the bam file that is to get merged into bam_file_1
- **bam_file_3** (*str*) – Location of the bam file that is to get merged into bam_file_1

bam_merge_4 (***kwargs*)

Wrapper for the pysam SAMtools merge function

Parameters

- **bam_file_1** (*str*) – Location of the bam file to merge into
- **bam_file_2** (*str*) – Location of the bam file that is to get merged into bam_file_1
- **bam_file_3** (*str*) – Location of the bam file that is to get merged into bam_file_1
- **bam_file_4** (*str*) – Location of the bam file that is to get merged into bam_file_1

bam_merge_5 (***kwargs*)

Wrapper for the pysam SAMtools merge function

Parameters

- **bam_file_1** (*str*) – Location of the bam file to merge into
- **bam_file_2** (*str*) – Location of the bam file that is to get merged into bam_file_1
- **bam_file_3** (*str*) – Location of the bam file that is to get merged into bam_file_1
- **bam_file_4** (*str*) – Location of the bam file that is to get merged into bam_file_1
- **bam_file_5** (*str*) – Location of the bam file that is to get merged into bam_file_1

bam_paired_reads (***kwargs*)

Wrapper for the pysam SAMtools view function to identify if a bam file contains paired end reads

Parameters `bam_file` (*str*) – Location of the bam file that is to be indexed

Returns True if the bam file contains paired end reads

Return type bool

bam_sort (***kwargs*)

Wrapper for the pysam SAMtools sort function

Parameters `bam_file` (*str*) – Location of the bam file to sort

bam_stats (***kwargs*)

Wrapper for the pysam SAMtools flagstat function

Parameters

- `bam_file` (*str*) – Location of the bam file that is to be indexed
- `bam_idx_file` (*str*) – Location of the bam index file (.bai)

check_header (***kwargs*)

Wrapper for the pysam SAMtools merge function

Parameters

- `bam_file_1` (*str*) – Location of the bam file to merge into
- `bam_file_2` (*str*) – Location of the bam file that is to get merged into `bam_file_1`

5.4 FASTQ Functions

The following functions are ones that are used for the manipulation of FASTQ files.

5.4.1 Reading

The following functions are to provide easy access for iterating through entries within a FASTQ file(s) both single and paired.

class `tool.fastqreader.fastqreader`

Module for reading single end and paired end FASTQ files

closeFastQ ()

Close file handles for the FastQ files.

closeOutputFiles ()

Close the output file handles

createOutputFiles (*tag=""*)

Create and open the file handles for the output files

Parameters `tag` (*str*) – Tag to identify the output files (DEFAULT: ‘')

eof (*side=1*)

Indicate if the end of the file has been reached

Parameters `side` (*int*) – 1 or 2

incrementOutputFiles ()

Increment the counter and create new files for splitting the original FastQ paired end files.

next (*side=1*)

Get the next read element for the specific FastQ file pair

Parameters `side` (*int*) – 1 or 2 to get the element from the relevant end (DEFAULT: 1)

Returns

`id` [str] Sequence ID
`seq` [str] Called sequence
`add` [str] Plus sign
`score` [str] Base call score

Return type dict

openFastQ (*file1, file2=None*)

Create file handles for reading the FastQ files

Parameters

- `file1` (*str*) – Location of the first FASTQ file
- `file2` (*str*) – Location of a paired end FASTQ file.

writeOutput (*read, side=1*)

Writer to print the extracted lines

Parameters

- `read` (*dict*) – Read is the dictionary object returned from `self.next()`
- `side` (*int*) – The side that the read has come from (DEFAULT: 1)

Returns False if a value other than 1 or 2 is entered for the side.

Return type bool

5.4.2 Splitting

This tool has been created to aid in splitting FASTQ files into manageable chunks for parallel processing. It is able to work on single and paired end files.

class `tool.fastq_splitter.fastq_splitter` (*configuration=None*)

Script for splitting up FASTQ files into manageable chunks

paired_splitter (***kwargs*)

Function to divide the paired-end FastQ files into separate sub files of 1000000 sequences so that the aligner can run in parallel.

Parameters

- `in_file1` (*str*) – Location of first paired end FASTQ file
- `in_file2` (*str*) – Location of second paired end FASTQ file
- `tag` (*str*) – DEFAULT = tmp Tag used to identify the files. Useful if this is getting run manually on a single machine multiple times to prevent collisions of file names

Returns

- **Returns** (*Returns a list of lists of the files that have been generated.*) – Each sub list containing the two paired end files for that subset.
- `paired_files` (*list*) – List of lists of pair end files. Each sub list containing the two paired end files for that subset.

run (*input_files*, *input_metadata*, *output_files*)

The main function to run the splitting of FASTQ files (single or paired) so that they can be aligned in a distributed manner

Parameters

- **input_files** (*dict*) – List of input fastq file locations
- **metadata** (*dict*) –
- **output_files** (*dict*) –

Returns

- **output_file** (*str*) – Location of compressed (.tar.gz) of the split FASTQ files
- **output_names** (*list*) – List of file names in the compressed file

single_splitter (***kwargs*)

Function to divide the FastQ files into separate sub files of 1000000 sequences so that the aligner can run in parallel.

Parameters

- **in_file1** (*str*) – Location of first FASTQ file
- **tag** (*str*) – DEFAULT = tmp Tag used to identify the files. Useful if this is getting run manually on a single machine multiple times to prevent collisions of file names

Returns

- **Returns** (*Returns a list of the files that have been generated.*) – Each sub list containing the two paired end files for that subset.
- **paired_files** (*list*) – List of lists of pair end files. Each sub list containing the two paired end files for that subset.

5.4.3 Entry Functions

The following functions allow for manipulating FASTQ files.

class `tool.fastq_utils.fastqUtils`

Set of methods to help with the management of FastQ files.

static `fastq_match_paired_ends` (*fastq_1*, *fastq_2*)

Take 2 fastq files and remove ends that don't have a matching pair. Requires that the FastQ files are ordered correctly.

Mismatches can occur if there is a filtering step that removes one of the paired ends

Parameters

- **fastq_1** (*str*) – Location of FastQ file
- **fastq_2** (*str*) – Location of FastQ file

static `fastq_randomise` (*fastq*, *output=None*)

Randomising the order of reads within a FastQ file

Parameters

- **fastq** (*str*) – Location of the FastQ file to randomise
- **output** (*str*) – [OPTIONAL] Location of the output FastQ file. If left blank then the randomised output is saved to the same location as *fastq*

static fastq_sort_file (*fastq*, *output=None*)

Sorting of a FastQ file

Parameters

- **fastq** (*str*) – Location of the FastQ file to sort
- **output** (*str*) – [OPTIONAL] Location of the output FastQ file. If left blank then the sorted output is saved to the same location as *fastq*

Continuous Integration with Travis

This document summarizes the steps involved in setting up a continuous integration suite with Travis for our pipelines.

6.1 Getting Started

Login with your gitHub account details on [travis.ci](https://travis-ci.com). Add the “Multiscale Genomics” to your organization to gain access to its repositories. (you would have to send a request to be added to this organization).

Follow the onscreen instructions on [travis.ci](https://travis-ci.com).

```
1 Flick the repository switch to "on".
2 Add .travis.yml to your repository.
3 Sync your travis account to your GitHub account.
4 Pushing to Git will trigger the first Travis build after adding the above file.
```

6.2 Making .travis.yml File

To your added .travis.yml file in your GitHub repository include:

- “python” in “language”. With version/s specified in “python: “
- All packages required for running the pipelines in “addons: apt; packages: “. (For more information on the packages please see Full Installation from : http://multiscale-genomics.readthedocs.io/projects/mg-process-fastq/en/latest/full_installation.html#setup-the-system-environment)
- Docker in “services” to tell Travis it needs to have docker installed.
- **services:**
 - docker

- All tools to be installed within “install:” section (For more information on the packages please see Full Installation from : http://multiscale-genomics.readthedocs.io/projects/mg-process-fastq/en/latest/full_installation.html#setup-the-system-environment)

Note: libtbb did not seem to be installing correctly when put in “apt: packages: “. It has therefore been done with sudo in “install:”

- Setup all symlinks in “before_script: “
- Change execution permissions on shims folder and harness.sh
- Add the shims folder path to your \$PATH
- Include harness.sh to your “script” to be executed

6.3 Making harness.sh File

Include all test scripts to be tested with pytest to your harness.sh file. As iNPS does not work with Python < 3, a conditional check is present to ensure that it does not run unless Travis is running Python 3

6.4 Running Docker container

To run docker within Travis, it has been included within the “services” in the .yml file. Every time travis runs docker it will pull the latest publicly available image from DockerHub and run it. This gets done from within the shims files. The pre-built container can also be pulled from : <https://hub.docker.com/r/multiscalegenomics/mgprocessfastq/> using command :

```
docker pull multiscalegenomics/mgprocessfastq:testdocker
```

For more details on the docker container, please refer to : <https://github.com/Multiscale-Genomics/mg-process-fastq/blob/master/docs/docker.rst>

6.5 Setting up Shims

Libmaus2 and Biobambam2 have had to be installed within a docker container, as they were causing Travis to time out. As the container is non-interactable while on Travis and is not hosting live server. An intermediate layer to access the contents within docker from travis has been introduced in the form of shims. To make these files, take the list of all biobambam2 modules and construct bash script files with the following command for each of the modules :

```
exec docker run -it multiscalegenomics/mgprocessfastq:biobambamimage ~/lib/  
↳biobambam2/bin/biobambam_module_name $@
```

The .travis.yml file used for testing the mg-process-fastq pipelines can be found at : <https://github.com/Multiscale-Genomics/mg-process-fastq/blob/master/.travis.yml>

Setting up and using a Docker Container

7.1 Our reason for using a container

While working with Travis, the installation of libmaus2 and biobambam2 took up more than 45 minutes (accumulative with the rest of the installations), which caused Travis to time out. We therefore resorted to putting both tools in a container and accessing the commands from there. This document summarizes the steps involved in making a docker image for the above two tools and running a container from that image. As well as uploading your image to docker hub to make it publicly accessible.

This document has been prepared keeping macOS Sierra in mind, although many of the commands are cross platform (*nix) compliant.

7.2 Getting Started

To be able to build a docker container you must have :

- a) Docker installed on your machine
- b) An account on one of the docker repositories (Docker Hub or Quay). We have used Docker Hub as this was free access.

7.2.1 a) Installing docker to your machine

For this work I had installed a command line based docker, along with the Virtual machine boot2docker. There is however a GUI distribution available for MAC as well. You may install boot2docker using :

```
1 brew install boot2docker
```

7.2.2 b) Setting up account on Docker Hub

Go to <https://hub.docker.com> and setup an account with Docker Hub. Add Multiscale Genomics to your organizations and create a repository : mgprocessfastq. You will be uploading your docker images to this repository later.

7.3 Constructing a docker container

Run the following preliminary commands to get your boot2docker running:

```
1 boot2docker up
2 eval "$ (boot2docker shellinit) "
```

You would also need to have :

- docker-machine
- Virtual Box

installed on your Mac. After these, execute the following commands:

```
1 docker-machine create -d virtualbox dev
2 eval $(docker-machine env dev)
```

To ensure your docker is running:

```
1 docker
```

7.3.1 Making the Dockerfile for libmaus2 and biobambam2

You may want to create a new folder for this purpose, as the docker command compiles the Dockerfile with the given path to the folder. Create a new file with the name of “Dockerfile”. Include the following lines within this file:

```
FROM ubuntu:14.04

RUN apt-get update && \
    apt-get -y install sudo

RUN sudo apt-get install -y make build-essential libssl-dev zlib1g-dev \
libbz2-dev libreadline-dev libsqlite3-dev wget curl llvm libncurses5-dev \
libncursesw5-dev xz-utils tk-dev unzip mcl libgtk2.0-dev r-base-core \
libcurl4-gnutls-dev python-rpy2 git

RUN mkdir Mug \
    && cd Mug \
    && apt-get -y install git \
    && git config --global user.name "your_username" \
    && git config --global user.email "your_emailId" \
    && pwd \
    && mkdir bin lib code \
    && cd lib \
    && git clone https://github.com/gt1/libmaus2.git
    && cd libmaus2 \
    && sudo apt-get -y install libtool m4 automake \
    && libtoolize \
    && aclocal \
```

(continues on next page)

(continued from previous page)

```
&& autoheader \
&& automake --force-missing --add-missing \
&& autoconf \
&& ./configure --prefix=/Mug/lib/libmaus2 \

&& make \
&& make install \
&& cd /Mug/lib \

&& git clone https://github.com/gt1/biobambam2.git && cd biobambam2 \
&& autoreconf -i -f \
&& ./configure --with-libmaus2=/Mug/lib/libmaus2 --prefix=/Mug/lib/biobambam2 \
&& make install
```

7.3.2 Making the docker image

Build a docker image from this file using:

```
1 cd /path/to/your/dockerfile
2 docker build -t multiscalegenomics/mgprocessfastq/biobambamimage.
```

Login with your docker hub account details :

```
1 docker login
```

Push the above image to your docker hub repository

```
1 docker push multiscalegenomics/mgprocessfastq:biobambamimage
```

7.3.3 Running a docker container

You should be able to run the above image locally on your machine as well as pulling it elsewhere (on a system which has docker):

```
1 docker pull multiscalegenomics/mgprocessfastq:biobambamimage
```

and then running a container via :

```
1 docker run --name name_you_want multiscalegenomics/mgprocessfastq:biobambamimage
```

Our Travis build pulls the image from our mgprocessfastq repository from within the shims files, and runs the containers using the commands within.

8.1 2017-08-15 - Implementation of pigz

It was highlighted that the archival step using the python module tarfile was taking a long time when archiving and compressing large volumes of data. This is an issue in the WGBS pipeline where the FASTQ files are broken down into small sections and are then aligned individually.

The issue was in relation to the compression step. tar, and therefore tarfile, runs as a single process. In linear time to compress a 14GB file takes 18 minutes, which is a long time for the user that have to wait. The archival step is less than 1 min. As a result the compression has moved to using pigz after the archival step to perform the compression in a parallel fashion across all available cores on a machine. On a 4 core machine this allowed the compression of the same file to take only 7 min.

The resultant compressed file remains accessible via gzip. After decompression files have the same structure, so this should be an invisible change.

8.2 2018-01-26 - Disable no-self-use for @tasks

The disabling of pylint test for determining if a function should be a static method was made as this affected the behaviour of pycompss preventing it from functioning correctly. As a result all @task functions do not require the no-self-use test to be run.

8.3 2018-02-28 - BAM Merge Strategy

Based on benchmarking of the pipeline the procedure for merging of bam files has been modified to get an optimal balance between running as much as possible in parallel vs the cost of spinning up new jobs to perform the merging. It was found that running each job merging 10 files provided the break even point between the cost of creating a new job and getting the maximum throughput through the system. It also reduces the number of iterative merge procedures which is beneficial when there are alignments that are difficult to merge.

8.4 2018-04-26 - BAM Splitting

Added the required functions for tools to be able to split a bam by the number of chromosomes so that the analysis can be done in parallel. As an initial run this has been implemented in the MACS2 tool, where the indexing of the bam file is performed by the tool. The bam and bai files are passed to the jobs so that they can then extract the chromosome that they are required to analyse.

In the future the creation of the bai file could be done at alignment time, but in the case of MACS2 there is a filtering step on the aligned bam file, so a new index would be required.

8.5 2018-05-01 - Compression of FASTQ

Added compression of the split FASTQ files to reduce the amount of space required when processing the data. There is also the removal of the tmp directory after the completion of the splitter to save on space.

8.6 2018-05-09 - Handling aligner index decompression

The code has been modified so that there is a single decompression of the BWA and Bowtie2 common indexes. The index files are then explicitly handed to the alignment task rather than handing over the compressed index. The decompression is performed as a @task so that the index files are already in the COMPSs system. This means that handing the index files to the alignment tasks creates a single symlink in the sandbox temporary file directory rather than duplicating the whole of the index structure for each job.

8.7 2018-05-22 - GEM Naming

Update so that the gem files are name <genome-file>.gem.gz inline with requests from WP7 partners so that the name of the index is picked up correctly

8.8 2018-05-22 - TrimGalore

To try and improve the quality of the reads that are used for numerous pipelines, TrimGalore has been included as a pipeline to aid in the clipping and removal of low quality regions of reads. The pipeline can be run on single or paired end FASTQ files. A report of the trimmed data is also returned for the user to identify what changes were made.

8.9 2018-05-31 - Public genomes and indexes

The VRE is making it possible to use centrally stored genomic sequences and the pre-built indexes. This relies on a second key-value pair for the genome and indexes. Pipelines that interact with the public files need to be able to test for the present of “genome_public” and “index_public”, if they are not present then they should be able to safely fall back on “genome” and “index” keys-value pairs. With the separation of the tools and the pipelines this means that the changes should only need to happen in the pipelines as this can perform the translation between the “genome_public” to “genome” ready for the tool.

8.10 2018-06-01 - Separated WGBS Code Testing

To bring down the run time for the TravisCI, the WGBS has been moved to a separate track. This has the benefit of getting the testing started earlier and allowing the other tests to finish sooner.

8.11 2018-06-01 - Travis Caching

Travis CI is now able to cache the lib directory so that the build phase is reduced to improve test speeds. If there are changes to the lib then these need to be refreshed in the TravisCI settings to ensure that the new libraries are included, or flushed if there are changes to the versions of packages in the lib.

There is also caching of the pip directory to reduce the load time.

8.12 2018-06-04 - Split the WGBS test scripts

Split the testing of the WGBS pipeline and tool chains so that they 2 sets can run in parallel. Both take too long when run in series.

8.13 2018-06-05 - Use of the logger PROGRESS

Added in the use of the logger.progress to indicate the progression of a process.

8.14 2018-06-14 - Paired end alignment

The aligner pipelines has been modified the pass through all the input and metadata to the aligner tools, this simplifies the the passing of a second fastq file and also make using these pipelines for alignment of paired end data possible.

8.15 2018-06-18 - Branch tidying during alignment

Modified the way that the alignment pipelines manage the temporary files. These are now deleted once the pipeline has finished using them. The purpose of this is to save space on the file system and prevent large jobs taking up too much space.

There have also been changes to the handling of paired end files for the alignment pipelines improving the clarity of what is happening and simplifying the passing of parameters. There are also changes to the tests to allow for the removal of temporary files and there are tests to make sure that the output bam files are single or paired end.

Other changes include: - Simplification of the untarring functions - Modifications to the Bowtie2 index file for consistency with the BWA index file - Refactored the BWA ALN sai file generation to reduce redundancy to allow for multi-processing when there is paired-end data - Improved the handling of the suffixes for FASTQ and FASTA files so that it can handle variants

8.16 2018-06-27 - Remove reads marked as duplicate by BioBamBam

BioBamBam only marks reads as duplicate, but does not remove the after. The Tool has been updated to remove the flagged duplicates using samtools with the parameter *-F 1024*. This matches the pipeline used within the [Blueprints project](#).

Also performed some tidying of the code to annotate issues that had been highlighted by pylint.

8.17 2018-07-11 - Changes FASTQ splitter file management

The previous splitter would split the FASTQ files into separate changes, then create the tar file and then the Gzip file. This results in a large amount of wasted tmp space, which is a limited resource. The changes implemented incrementally add the sub-FASTQ files to the archive file, deleting them once they have been added. The whole archive file is then compressed. This has a large advantage when handling larger human datasets.

There has also been some refactoring of the handling of the archiving and compression steps to reduce the duplication of code within the repository.

8.18 2018-07-16 - Modified handling of file locations

Updated the handling of file locations to use `os.path.join` and `os.path.split` to allow for compatibility between different operating systems for the pipelines and tools.

8.19 2018-08-02 - Added in Paired End BAM file handling for MACS2

MACS2 is able to automatically handle the files that are handed to it except for paired-end BAM and BED files (BAMPE and BEDPE respectively). The MACS2 tool only accepts BAM files so a check was implemented to determine if the BAM file contained paired-end reads.

There has also been a major rewrite of the MACS2 tool to remove code duplication.

8.20 2018-07-16 - Modified handling of file locations

Updated the handling of file locations to use `os.path.join` and `os.path.split` to allow for compatibility between different operating systems for the pipelines and tools.

8.21 2018-08-07 - Storing tool parameters as part of the metadata

To improve the amount of information that is stored about the run of a tool, the parameters that were used are now being included as part of the metadata.

8.22 2018-08-07 - Extra output files from MACS2

MACS2 is able to generate a plot of the results as well as a bedGraph. These have now been integrated as part of the output files from the tool.

8.23 2018-08-13 - Normalised the use of OSError

IOError was deprecated in favour of OSError when moving to py3, but to maintain backwards compatibility IOError also needs to be supported. There were places in the code where this was not true and other places that relied on just OSError. Instances of just IOError have been converted to testing for both IOError and OSError and visa versa.

8.24 2018-08-15 - Use the config.json execution path

Using the directory of the input file for building the location of the working directory with outside of a task is not a viable option as it can write data to the wrong directory. The execution path provided in the config.json arguments is the right place. This location is also the location for output files. This issue occurred as the FASTQ splitter was generating a tar file that the aligners were downloading to the wrong location. Even though this was tidied up this was still not the right place to put this file.

8.25 2018-08-16 - Prevent further duplicate filtering by MACS2

In the process_chipseq.py pipeline the duplicates have already been filtered by BioBamBam2 and samtools so there is no need for further filtering to be done by MACS2.

8.26 2018-09-04 - Adding functionality to bam_utils and MACS2

MACS2 was previously set to work with the BAMPE option for the `-f/--format` parameter. Additional functionality has been added to `bam_utils` and `macs2` mode to incorporate the BEDPE option. This has been done for the Atac Seq pipeline to incorporate the processing of bed file rather than bam files if the user would need changes to the result files generated.

8.27 2018-09-17 - Updates to tool and pipeline run()

Changes to the pipelines so that the `run()` function matches the definitions within the Tool API. There have also been a number of changes so that the pipeline and tool code is python 3 compatible

8.28 2018-08-22 - Improvement of tadbit tools wrappers

A json with the matrix was included in the outputs of tadbit bin New normalization method OneD in tadbit normalize Code update to use last features of the development branch of tadbit tools api The wrapper of tadbit model was rebuilt to allow the modelling of full genomes, mainly for yeast General reshape of all the code according to pylint Inclusion of tests for the wrappers and tools of the tadbit pipelines

8.29 2018-09-25 - Converting the Kallisto TSV file to BED

To display the scores on the genome browser the abundance tsv is used to generate a bed file where the score matches the transcripts per million column from the abundance.tsv output from Kallisto. This module requires the presence of the `ensembl gff3` file for the matching assembly. This should be passed by the VRE when passing the FASTA file for the transcripts.

8.30 2018-10-18 - Multi File handling for the DamID-seq Pipeline

Ability to handle multiple input single/paired end data and background data files and process them in an orderly fashion. Ability to handle the resultant multiples of generated bam files in the idear tool and its matching individual pipeline.

8.31 2018-10-25 - WGBS Pipeline Create BigWig files as standard

The output wig files from BS Seeker2 are now converted to BigWig files by default rather than returning wig files. This is so that they are easier to visualise on the JBrowse interface.

8.32 2018-10-31 - Modify the file names for docker script

There were inconsistencies in the names, these have been updated to be more consistent with the consortium naming instead.

8.33 2018-11-08 - Modifications for the movement of files

To keep the memory usage low the files are now moved in N byte chunks rather than as a whole file.

8.34 2018-11-16 - Reading of Gzipped FASTQ files

Due to the limitations of space it is necessary for the users to be able to load only the gzipped versions of the FASTQ files. Added the ability to read the gzipped FASTQ files to the FASTQ reader. This is behind the aligners and the splitting procedure so has a wide ranging benefit for pipelines.

Apache License Version 2.0, January 2004 <http://www.apache.org/licenses/>

1. Definitions.

“License” shall mean the terms and conditions for use, reproduction, and distribution as defined by Sections 1 through 9 of this document.

“Licensor” shall mean the copyright owner or entity authorized by the copyright owner that is granting the License.

“Legal Entity” shall mean the union of the acting entity and all other entities that control, are controlled by, or are under common control with that entity. For the purposes of this definition, “control” means (i) the power, direct or indirect, to cause the direction or management of such entity, whether by contract or otherwise, or (ii) ownership of fifty percent (50%) or more of the outstanding shares, or (iii) beneficial ownership of such entity.

“You” (or “Your”) shall mean an individual or Legal Entity exercising permissions granted by this License.

“Source” form shall mean the preferred form for making modifications, including but not limited to software source code, documentation source, and configuration files.

“Object” form shall mean any form resulting from mechanical transformation or translation of a Source form, including but not limited to compiled object code, generated documentation, and conversions to other media types.

“Work” shall mean the work of authorship, whether in Source or Object form, made available under the License, as indicated by a copyright notice that is included in or attached to the work (an example is provided in the Appendix below).

“Derivative Works” shall mean any work, whether in Source or Object form, that is based on (or derived from) the Work and for which the editorial revisions, annotations, elaborations, or other modifications represent, as a whole, an original work of authorship. For the purposes of this License, Derivative Works shall not include works that remain separable from, or merely link (or bind by name) to the interfaces of, the Work and Derivative Works thereof.

“Contribution” shall mean any work of authorship, including the original version of the Work and any modifications or additions to that Work or Derivative Works thereof, that is intentionally submitted to Licensor for inclusion in the Work by the copyright owner or by an individual or Legal Entity authorized to submit on behalf of the copyright owner. For the purposes of this definition, “submitted” means any form of electronic, verbal, or written communication sent to the Licensor or its representatives, including but not limited to communication on electronic mailing lists, source code control systems, and issue tracking systems that are managed by, or on behalf of, the Licensor for the purpose of discussing and improving the Work, but excluding communication that is conspicuously marked or otherwise designated in writing by the copyright owner as “Not a Contribution.”

“Contributor” shall mean Licensor and any individual or Legal Entity on behalf of whom a Contribution has been received by Licensor and subsequently incorporated within the Work.

2. **Grant of Copyright License.** Subject to the terms and conditions of this License, each Contributor hereby grants to You a perpetual, worldwide, non-exclusive, no-charge, royalty-free, irrevocable copyright license to reproduce, prepare Derivative Works of, publicly display, publicly perform, sublicense, and distribute the Work and such Derivative Works in Source or Object form.
3. **Grant of Patent License.** Subject to the terms and conditions of this License, each Contributor hereby grants to You a perpetual, worldwide, non-exclusive, no-charge, royalty-free, irrevocable (except as stated in this section) patent license to make, have made, use, offer to sell, sell, import, and otherwise transfer the Work, where such license applies only to those patent claims licensable by such Contributor that are necessarily infringed by their Contribution(s) alone or by combination of their Contribution(s) with the Work to which such Contribution(s) was submitted. If You institute patent litigation against any entity (including a cross-claim or counterclaim in a lawsuit) alleging that the Work or a Contribution incorporated within the Work constitutes direct or contributory patent infringement, then any patent licenses granted to You under this License for that Work shall terminate as of the date such litigation is filed.
4. **Redistribution.** You may reproduce and distribute copies of the Work or Derivative Works thereof in any medium, with or without modifications, and in Source or Object form, provided that You meet the following conditions:
 - (a) You must give any other recipients of the Work or Derivative Works a copy of this License; and
 - (b) You must cause any modified files to carry prominent notices stating that You changed the files; and
 - (c) You must retain, in the Source form of any Derivative Works that You distribute, all copyright, patent, trademark, and attribution notices from the Source form of the Work, excluding those notices that do not pertain to any part of the Derivative Works; and
 - (d) If the Work includes a “NOTICE” text file as part of its distribution, then any Derivative Works that You distribute must include a readable copy of the attribution notices contained within such NOTICE file, excluding those notices that do not pertain to any part of the Derivative Works, in at least one of the following places: within a NOTICE text file distributed as part of the Derivative Works; within the Source form or documentation, if provided along with the Derivative Works; or, within a display generated by the Derivative Works, if and wherever such third-party notices normally appear. The contents of the NOTICE file are for informational purposes only and do not modify the License. You may add Your own attribution notices within Derivative Works that You distribute, alongside or as an addendum to the NOTICE text from the Work, provided that such additional attribution notices cannot be construed as modifying the License.

You may add Your own copyright statement to Your modifications and may provide additional or different license terms and conditions for use, reproduction, or distribution of Your modifications, or for any such Derivative Works as a whole, provided Your use, reproduction, and distribution of the Work otherwise complies with the conditions stated in this License.

5. Submission of Contributions. Unless You explicitly state otherwise, any Contribution intentionally submitted for inclusion in the Work by You to the Licensor shall be under the terms and conditions of this License, without any additional terms or conditions. Notwithstanding the above, nothing herein shall supersede or modify the terms of any separate license agreement you may have executed with Licensor regarding such Contributions.
6. Trademarks. This License does not grant permission to use the trade names, trademarks, service marks, or product names of the Licensor, except as required for reasonable and customary use in describing the origin of the Work and reproducing the content of the NOTICE file.
7. Disclaimer of Warranty. Unless required by applicable law or agreed to in writing, Licensor provides the Work (and each Contributor provides its Contributions) on an “AS IS” BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied, including, without limitation, any warranties or conditions of TITLE, NON-INFRINGEMENT, MERCHANTABILITY, or FITNESS FOR A PARTICULAR PURPOSE. You are solely responsible for determining the appropriateness of using or redistributing the Work and assume any risks associated with Your exercise of permissions under this License.
8. Limitation of Liability. In no event and under no legal theory, whether in tort (including negligence), contract, or otherwise, unless required by applicable law (such as deliberate and grossly negligent acts) or agreed to in writing, shall any Contributor be liable to You for damages, including any direct, indirect, special, incidental, or consequential damages of any character arising as a result of this License or out of the use or inability to use the Work (including but not limited to damages for loss of goodwill, work stoppage, computer failure or malfunction, or any and all other commercial damages or losses), even if such Contributor has been advised of the possibility of such damages.
9. Accepting Warranty or Additional Liability. While redistributing the Work or Derivative Works thereof, You may choose to offer, and charge a fee for, acceptance of support, warranty, indemnity, or other liability obligations and/or rights consistent with this License. However, in accepting such obligations, You may act only on Your own behalf and on Your sole responsibility, not on behalf of any other Contributor, and only if You agree to indemnify, defend, and hold each Contributor harmless for any liability incurred by, or claims asserted against, such Contributor by reason of your accepting any such warranty or additional liability.

END OF TERMS AND CONDITIONS

APPENDIX: How to apply the Apache License to your work.

To apply the Apache License to your work, attach the following boilerplate notice, with the fields enclosed by brackets “{}” replaced with your own identifying information. (Don’t include the brackets!) The text should be enclosed in the appropriate comment syntax for the file format. We also recommend that a file or class name and description of purpose be included on the same “printed page” as the copyright notice for easier identification within third-party archives.

Copyright [yyyy] [name of copyright owner]

Licensed under the Apache License, Version 2.0 (the “License”); you may not use this file except in compliance with the License. You may obtain a copy of the License at

<http://www.apache.org/licenses/LICENSE-2.0>

Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an “AS IS” BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

To ensure that the code written to process the data within the MuG VRE works it is essential to have a testing framework to ensure that as the code evolves it does not break the functionality of the pipelines and tools. There are 2 key parts to this:

1. Sample data
2. Runnable scripts

For each tool within the mg-process-fastq repository there is an initial set of sample data and matching tests for each tool and pipeline.

10.1 Sample Data

10.1.1 Test Data for ChIP-seq pipeline

The following document is for the preparation of data set required for testing the ChIP-seq pipeline. The document has been written with macOS Sierra in mind, although many of the commands are cross platform (*nix) compliant.

You would need to have the tools listed in “Prerequisites” installed on your system. For more details on installing the tools for this pipeline please refer to

http://multiscale-genomics.readthedocs.io/projects/mg-process-fastq/en/latest/full_installation.html

If you already have certain packages installed feel free to skip over certain steps. Likewise the bin, lib and code directories are relative to the home dir; if this is not the case for your system then make the required changes when running these commands.

Prerequisites

- BWA
- MACS 2
- Biobambam

- Samtools

Data set for genome file

Filtering for required coverage

Download the genome file from

```
wget "http://www.ebi.ac.uk/ena/data/view/CM000663.2,CM000664.2,CM000665.2,CM000666.2,
↪CM000667.2,CM000668.2,CM000669.2,CM000670.2,CM000671.2,CM000672.2,CM000673.2,
↪CM000674.2,CM000675.2,CM000676.2,CM000677.2,CM000678.2,CM000679.2,CM000680.2,
↪CM000681.2,CM000682.2,CM000683.2,CM000684.2,CM000685.2,CM000686.2,J01415.2&
↪display=fasta&download=fasta&filename=entry.fasta" -O GCA_000001405.22.fasta
```

Checkout https://github.com/Multiscale-Genomics/mg-misc-scripts/blob/master/ChIPSeq_Scripts/extract_chromosomeForChIP.py and extract chromosome 22 from the above file using the following command.

```
python extract_chromosomeForChIP.py path/to/your/input/file path/to/output/file
```

Download the fastq file from

```
wget ftp://ftp.sra.ebi.ac.uk/vol1/fastq/DRR000/DRR000150/DRR000150.fastq.gz
```

Unzip this file.

```
gunzip DRR000150.fastq.gz
```

Index the fasta file

```
bwa index GCA_000001405.22.chr22.fa.fasta
```

Align the fastq file

```
bwa aln GCA_000001405.22.chr22.fa.fasta DRR000150.chr22.fastq >GCA_000001405.22.chr22.
↪sai
```

And make the sam file

```
bwa samse GCA_000001405.22.chr22.fa.fasta GCA_000001405.22.chr22.sai DRR000150.chr22.
↪fastq >GCA_000001405.22.chr22.sam
```

Sort the sam file

```
samtools sort GCA_000001405.22.chr22.sam >GCA_000001405.22.chr22.sorted.sam
```

Find the depths of coverage from the sorted file

```
samtools depth GCA_000001405.22.chr22.sorted.sam >GCA_000001405.22.chr22.dp
```

From the depth file, find regions with ≥ 70 depth, spanning over ≥ 55 base pairs. You may get the script for this from https://github.com/Multiscale-Genomics/mg-misc-scripts/blob/master/ChIPSeq_Scripts/traverseForCoverageRegion_ChIP.py. Run it using:

```
python traverseForCoverageRegion_ChIP.py path/to/GCA_000001405.22.chr22.dp
```


Running this script would print the spanning regions. If it is a continuous region, you may only take the first starting base pair and the last ending base pair, as inputs for the next step. (Take out 1000 and add in 1000 to these respectively to get upstream and downstream spanning bases)

Extract the corresponding fasta sequence from the chromosome file for the positions retrieved from the above step. Checkout file from https://github.com/Multiscale-Genomics/mg-misc-scripts/blob/master/ChIPSeq_Scripts/extractChromosomalRegion.py and run using command:

```
python extractChromosomalRegion.py path/to/original/fasta/file path/to/output/file/  
↪for/region/macs2.Human.GCA_000001405.22.fasta starting_base_position ending_base_  
↪position
```

Making the Fastq file

Index the fasta file for the selected region

```
bwa index macs2.Human.GCA_000001405.22.fasta
```

Align the fastq file

```
bwa aln macs2.Human.GCA_000001405.22.fasta DRR000150.chr22.fastq >macs2.Human.GCA_  
↪000001405.22.sai
```

And make the sam file

```
bwa samse macs2.Human.GCA_000001405.22.fasta macs2.Human.GCA_000001405.22.sai_  
↪DRR000150.chr22.fastq >macs2.Human.GCA_000001405.22.sam
```

Filter this sam file for the reads which aligned with chromosome 22 using the following command:

```
awk '$3 == chr22' macs2.Human.GCA_000001405.22.sam >macs2.Human.GCA_000001405.22.22.  
↪sam
```

From the filtered reads from the above output file, extract the corresponding entries in fastq file. You may do this using the file at :

```
https://github.com/Multiscale-Genomics/mg-misc-scripts/blob/master/ChIPSeq\_Scripts/  
↪makeFastQFiles.py
```

and running it via command line:

```
python makeFastQFiles.py --samfile macs2.Human.GCA_000001405.22.22.sam --fastQfile /  
↪path/to/DRR000150.chr22.fastq --pathToOutput /path/to/save/output/fastq/file/to/ --  
↪fastqOut macs2.Human.DRR000150.22.fastq
```

The fastq file in the above step and fasta file macs2.Human.GCA_000001405.22.fasta together make up the data set for ChIP-seq pipeline

10.1.2 iDamID-Seq Test Data

Test Data

Dataset

Stable ID	SRR3714775
Citation	Franks et al 2016

Genome

Assembly	GRCh38
Chromosome	22
Start	15000000
End	19999999

Method

The full dataset was downloaded from ENA aligned to the genome using GEM.

```

1 wget ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR371/005/SRR3714775/SRR3714775.fastq.gz
2 wget ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR371/005/SRR3714775/SRR3714776.fastq.gz
3 wget ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR371/005/SRR3714775/SRR3714777.fastq.gz
4 wget ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR371/005/SRR3714775/SRR3714778.fastq.gz
5
6 wget "http://www.ebi.ac.uk/ena/data/view/CM000684.2&display=fasta&download=fasta&
   ↪filename=entry.fasta" -O GCA_000001405.22.chr22.fasta
7
8 bwa index GCA_000001405.22.chr22.fasta
9
10 bwa aln -q 5 -f SRR3714775.fastq.sai
11 bwa aln -q 5 -f SRR3714776.fastq.sai
12 bwa aln -q 5 -f SRR3714777.fastq.sai
13 bwa aln -q 5 -f SRR3714778.fastq.sai
14
15 bwa samse -f SRR3714775.fastq.sam GCA_000001405.22.chr22.fasta SRR3714775.fastq.sai_
   ↪SRR3714775.fastq
16 bwa samse -f SRR3714776.fastq.sam GCA_000001405.22.chr22.fasta SRR3714776.fastq.sai_
   ↪SRR3714776.fastq
17 bwa samse -f SRR3714777.fastq.sam GCA_000001405.22.chr22.fasta SRR3714777.fastq.sai_
   ↪SRR3714777.fastq
18 bwa samse -f SRR3714778.fastq.sam GCA_000001405.22.chr22.fasta SRR3714778.fastq.sai_
   ↪SRR3714778.fastq
19
20 samtools view -h -o SRR3714775.sam SRR3714775.bam
21 samtools view -h -o SRR3714776.sam SRR3714776.bam
22 samtools view -h -o SRR3714777.sam SRR3714777.bam
23 samtools view -h -o SRR3714778.sam SRR3714778.bam
24
25 awk 'BEGIN { FS = "[[:space:]]+" } ; $3 ~ /CM000684/ {print $1}' SRR3714775.sam >_
   ↪SRR3714775.chr22.sam
26 awk 'BEGIN { FS = "[[:space:]]+" } ; $3 ~ /CM000684/ {print $1}' SRR3714776.sam >_
   ↪SRR3714776.chr22.sam

```

(continues on next page)

(continued from previous page)

```

27 awk 'BEGIN { FS = "[[:space:]]+" } ; $3 ~ /CM000684/ {print $1}' SRR3714777.sam >_
    ↪SRR3714777.chr22.sam
28 awk 'BEGIN { FS = "[[:space:]]+" } ; $3 ~ /CM000684/ {print $1}' SRR3714778.sam >_
    ↪SRR3714778.chr22.sam
29
30 python scripts/ExtractRowsFromFASTQs.py --input_1 tests/data/SRR3714775.fastq --rows_
    ↪tests/data/SRR3714775.chr22.sam --output_tag profile
31 python scripts/ExtractRowsFromFASTQs.py --input_1 tests/data/SRR3714776.fastq --rows_
    ↪tests/data/SRR3714776.chr22.sam --output_tag profile
32 python scripts/ExtractRowsFromFASTQs.py --input_1 tests/data/SRR3714777.fastq --rows_
    ↪tests/data/SRR3714777.chr22.sam --output_tag profile
33 python scripts/ExtractRowsFromFASTQs.py --input_1 tests/data/SRR3714778.fastq --rows_
    ↪tests/data/SRR3714778.chr22.sam --output_tag profile

```

The alignments were then filtered with BioBamBam and peak calling was performed with iDEAR and a suitable region with a number of peaks was identified. The chromosomal region was extract and the matching reads to this region were identified. To reduce the number of reads that matched so that it could be used as a test set for the code base every other read was selected so that a reasonable number of reads we present. This mean that there are results when running the test, but generating the alignments does not take too long to compute.

```

1 sed -n 539916,556583p GCA_000001405.chr22.fasta > GCA_000001405.chr22ss.fasta
2 sed -n 5002,11669p idear.Human.GCA_000001405.22.fasta >> idear.Human.GCA_000001405.22.
    ↪subset.fasta

```

Test Scripts

The following are the tests for checking that the tools in the Hi-C pipeline are functioning correctly.

The tests should be run in this order so that the required input files are generated at the correct stage.

```

1 pytest -m idamidseq tests/test_bwa_indexer.py
2 pytest -m idamidseq tests/test_bwa_aligner.py
3 pytest -m idamidseq tests/test_biobambam.py
4 pytest -m idamidseq tests/test_bsgenome.py
5 pytest -m idamidseq tests/test_idear.py

```

These can be called as part of a single tool chain with:

```

1 python tests/test_toolchains.py --pipeline idamidseq

```

10.1.3 Test Data for MNase-seq pipeline

The following document is for the preparation of data set required for testing the MNase-seq pipeline. The document has been written with macOS Sierra in mind, although many of the commands are cross platform (*nix) compliant.

You would need to have the tools listed in “Prerequisites” installed on your system. For more details on installing the tools for this pipeline please refer to

http://multiscale-genomics.readthedocs.io/projects/mg-process-fastq/en/latest/full_installation.html

If you already have certain packages installed feel free to skip over certain steps. Likewise the bin, lib and code directories are relative to the home dir; if this is not the case for your system then make the required changes when running these commands.

Prerequisites

- BWA
- Samtools
- iNPS

Note: iNPS will only run within a Python 3 or above environment

Data set for genome file

Filtering for required coverage

Download the fasta file for Mouse chromosome 19 from

```
wget "http://www.ebi.ac.uk/ena/data/view/CM001012&display=fasta&download=fasta&
↪filename=CM001012.fasta" -O Mouse.CM001012.2.fasta
```

Download the fastq file from

```
wget ftp://ftp.sra.ebi.ac.uk/vol1/fastq/DRR000/DRR000386/DRR000386.fastq.gz
```

Unzip this file.

```
gunzip DRR000386.fastq.gz
```

Index the fasta file

```
bwa index Mouse.CM001012.2.fasta
```

Align the fastq file

```
bwa aln Mouse.CM001012.2.fasta DRR000386.fastq >Mouse.CM001012.2.sai
```

And make the sam file

```
bwa samse Mouse.CM001012.2.fasta Mouse.CM001012.2.sai DRR000386.fastq >Mouse.CM001012.
↪2.sam
```

Filter out the aligned reads from the above sam file.

```
awk '$3 != "*" ' Mouse.CM001012.2.sam >Mouse.CM001012.2.filtered.sam
```

Sort the sam file

```
samtools sort Mouse.CM001012.2.filtered.sam >Mouse.CM001012.2.sorted.sam
```

Find the depths of coverage from the sorted file

```
samtools depth Mouse.CM001012.2.sorted.sam >Mouse.CM001012.2.dp
```

From the depth file, find regions with ≥ 70 depth, spanning over ≥ 55 base pairs. You may get the script for this from: https://github.com/Multiscale-Genomics/mg-misc-scripts/blob/master/MNaseSeq_Scripts/traverseForCoverageRegion_MNase.py

Run it using:

```
python traverseForCoverageRegion_MNase.py path/to/Mouse.CM001012.2.dp
```

Running this script would print the spanning regions. Running this script for this data set gives multiple regions. The output is in the format : start - end - depth. The one at the end has a maximal coverage from this data set. Since it is a continuous region, you may take the first starting base pair and the last ending base pair, as inputs for the next step. (Take out 1000 and add in 1000 to these respectively to get upstream and downstream spanning bases)

Extract the corresponding fasta sequence from the chromosome file for the positions retrieved from the above step. Checkout file from https://github.com/Multiscale-Genomics/mg-misc-scripts/blob/master/MNaseSeq_Scripts/extractChromosomalRegion.py and run using command:

```
python extractChromosomalRegion.py path/to/original/fasta/file path/to/output/file/
↪for/region/inps.Mouse.GRCm38.fasta starting_base_position ending_base_position
```

Making the Fastq file

Index the fasta file for the selected region

```
bwa index inps.Mouse.GRCm38.fasta
```

Align the fastq file

```
bwa aln inps.Mouse.GRCm38.fasta DRR000386.fastq >inps.Mouse.GRCm38.sai
```

And make the sam file

```
bwa samse inps.Mouse.GRCm38.fasta inps.Mouse.GRCm38.sai DRR000386.fastq >inps.Mouse.
↪GRCm38.sam
```

Filter this sam file for the reads which aligned with chromosome 19 using the following command:

```
awk '$3 != "*" ' inps.Mouse.GRCm38.sam >inps.Mouse.GRCm38.sam.19.sam
```

From the filtered reads from the above output file, extract the corresponding entries in fastq file. You may do this using the file at https://github.com/Multiscale-Genomics/mg-misc-scripts/blob/master/MNaseSeq_Scripts/makeFastQFiles.py

and running it via command line:

```
python makeFastQFiles.py --samfile path/to/inps.Mouse.GRCm38.sam.19.sam --fastQfile /
↪path/to/DRR000386.fastq --pathToOutput /path/to/save/output/fastq/file/to/ --
↪fastqOut DRR000386.MNaseseq.fastq
```

Shorten this file by running the script at https://github.com/Multiscale-Genomics/mg-misc-scripts/blob/master/MNaseSeq_Scripts/randomSeqSelector.py

using

```
python randomSeqSelector.py DRR000386.MNaseseq.fastq inps.Mouse.DRR000386.fastq
```

The fastq file in the above step and fasta file inps.Mouse.GRCm38.fasta together make up the data set for MNase-seq pipeline

10.1.4 Test Data for RNA-seq pipeline

The following document is for the preparation of data set required for testing the RNA-seq pipeline. The document has been written with macOS Sierra in mind, although many of the commands are cross platform (*nix) compliant.

You would need to have the tools listed in “Prerequisites” installed on your system. For more details on installing the tools for this pipeline please refer to

http://multiscale-genomics.readthedocs.io/projects/mg-process-fastq/en/latest/full_installation.html

If you already have certain packages installed feel free to skip over certain steps. Likewise the bin, lib and code directories are relative to the home dir; if this is not the case for your system then make the required changes when running these commands.

Prerequisites

- Kallisto
- Samtools

Data set for genome file

Go to Ensemble website >> Human >> Example gene

```
http://www.ensembl.org/Homo_sapiens/Gene/Summary?g=ENSG00000139618;r=13:32315474-  
↪32400266
```

Copy the chromosome number and coordinates given in the “location” field. Go to BioMart (top panel), and select Filters from the left panel. Expand Regions and enter the information retrieved above.

Click on “Attributes” in the left panel. Select Gene stable ID, Transcript stable ID from Features. Select cDNA sequences from Sequences radio button.

Click on the Results button above the left panel. Export results to fasta file.

Index this file using Kallisto indexer:

```
kallisto index -i kallisto.Human.GRCh38.fasta.idx /path/to/file/exportSequences.fasta
```

Download the fastq files

```
wget ftp://ftp.sra.ebi.ac.uk/vol1/fastq/ERR030/ERR030872/ERR030872_1.fastq.gz  
wget ftp://ftp.sra.ebi.ac.uk/vol1/fastq/ERR030/ERR030872/ERR030872_2.fastq.gz
```

Run the Kallisto quantifier using command:

```
kallisto quant -i kallisto.Human.GRCh38.fasta.idx -o out --pseudobam /path/to/  
↪ERR030872_1.fastq.gz /path/to/ERR030872_2.fastq.gz >kallisto.ERR030872.sam
```

Filter the aligned sequence entries from the above sam file:

```
awk '$3 != "*" kallisto.ERR030872.sam >kallisto.ERR030872.filtered.sam
```

Unzip the fastq files.

```
unzip ERR030872_1.fastq.gz  
unzip ERR030872_2.fastq.gz
```

Checkout https://github.com/Multiscale-Genomics/mg-misc-scripts/blob/master/RNASeq_Scripts/makeFastQFiles.py and use the following command to generate the fastq files:

```
python /path/to/makeFastQFiles.py --samfile kallisto.ERR030872.filtered.sam --
↳fastQfile ERR030872_1.fastq --pathToOutput /path/to/make/fastqFile/ --fastqOut_
↳ERR030872_1.RNAseq.fastq
python /path/to/makeFastQFiles.py --samfile kallisto.ERR030872.filtered.sam --
↳fastQfile ERR030872_2.fastq --pathToOutput /path/to/make/fastqFile/ --fastqOut_
↳ERR030872_2.RNAseq.fastq
```

Shorten these files by running the script at https://github.com/Multiscale-Genomics/mg-misc-scripts/blob/master/RNASeq_Scripts/randomSeqSelector.py using

```
python PythonScripts/randomSeqSelector.py ERR030872_1.RNAseq.fastq kallisto.Human.
↳ERR030872_1.fastq
python PythonScripts/randomSeqSelector.py ERR030872_2.RNAseq.fastq kallisto.Human.
↳ERR030872_2.fastq
```

Then zip them:

```
gzip kallisto.Human.ERR030872_1.fastq
gzip kallisto.Human.ERR030872_2.fastq
```

10.1.5 WGBS Test Data

Test Data

Dataset

Stable ID	SRR892982
Citation	Sun et al 2014

Genome

Assembly	GRChm8, chr19
Chromosome	19
Start	3000000
End	4020000

Method

The full dataset was downloaded from ENA aligned to the genome using GEM.

```
1 wget ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR892/SRR892982/SRR892982_1.fastq.gz
2 wget ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR892/SRR892982/SRR892982_2.fastq.gz
3 gunzip SRR892982_1.fastq.gz
4 gunzip SRR892982_2.fastq.gz
5
6 wget "http://www.ebi.ac.uk/ena/data/view/CM001012&display=fasta&download=fasta&
↳filename=entry.fasta" -O mouse.GRCm38.19.fasta
```

(continues on next page)

(continued from previous page)

```

7
8 sed -n 40022200,41046060p GCA_000001635.6.fasta > mouse.GRCm38.19.fasta
9
10 bowtie2-build mouse.GRCm38.19.fasta mouse.GRCm38.19.idx
11 bowtie2 -x genome/mouse.GRCm38.19.idx -U data/SRR892982_1.fastq > SRR892982_1.sam
12 samtools view -h -o SRR892982_1.sam SRR892982_1.bam
13
14 awk 'BEGIN { FS = "[[:space:]]+" } ; $3 ~ /CM001012/ {print $1}' SRR892982_1.sam >
↳SRR892982_1.chr19.rows
15
16 python scripts/ExtractRowsFromFASTQs.py --input_1 tests/data/SRR892982_1.fastq --
↳input_2 tests/data/SRR892982_2.fastq --rows SRR892982_1.chr19.rows --output_tag
↳profile
17
18 mv data/tmp/SRR892982_1.profile_0.fastq data/SRR892982_1.chr19.fastq
19 mv data/tmp/SRR892982_2.profile_0.fastq data/SRR892982_2.chr19.fastq
20
21 bowtie2 -x genome/mouse.GRCm38.19.idx -1 data/SRR892982_1.chr19.fastq -2 data/
↳SRR892982_2.chr19.fastq > SRR892982_1.chr19.sam

```

There are a range of positive and negative peaks between 3000000 and 4020000. Subselect the genome and matching FASTQ reads for the required subsection:

```

1 head -n 1 mouse.GRCm38.19.fasta > 3M/mouse.GRCm38.19.3M.fasta
2 sed -n 50001,67001p mouse.GRCm38.19.fasta >> 3M/mouse.GRCm38.19.3M.fasta
3
4 bowtie2-build 3M/mouse.GRCm38.19.3M.fasta 3M/mouse.GRCm38.19.3M.idx
5 bowtie2 -p 4 -x 3M/mouse.GRCm38.19.3M.idx -1 SRR892982_1.chr19.fastq -2 SRR892982_2.
↳chr19.fastq > 3M/SRR892982.chr19.3M.sam
6 samtools view -h -o 3M/SRR892982.chr19.3M.sam 3M/SRR892982.chr19.3M.bam
7
8 python scripts/ExtractRowsFromFASTQs.py --input_1 SRR892982_1.chr19.fastq --input_2
↳SRR892982_2.chr19.fastq --rows 3M/SRR892982.chr19.3M.rows --output_tag subset
9
10 mv tmp/SRR892982_1.chr19.subset_0.fastq 3M/SRR892982_1.chr19.3M.fastq
11 mv tmp/SRR892982_2.chr19.subset_0.fastq 3M/SRR892982_2.chr19.3M.fastq

```

This has enough load to test the system, while also generating results in the output files.

These files are then saved in the tests/data directory as:

```

1 bsSeeker.Mouse.GRCm38.fasta
2 bsSeeker.Mouse.SRR892982_1.fastq.gz
3 bsSeeker.Mouse.SRR892982_2.fastq.gz

```

These files were too large for use within the Travis tst environment, so the number of entries was reduced by taking every other read:

```

1 sed -n -e '0~9{N;N;N;p}' tests/data/bsSeeker.Mouse.SRR892982_1.fastq > tests/data/
↳test_bsSeeker.Mouse.SRR892982_1.fastq
2 sed -n -e '0~9{N;N;N;p}' tests/data/bsSeeker.Mouse.SRR892982_2.fastq > tests/data/
↳test_bsSeeker.Mouse.SRR892982_2.fastq
3
4 mv tests/data/test_bsSeeker.Mouse.SRR892982_1.fastq tests/data/bsSeeker.Mouse.
↳SRR892982_1.fastq
5 mv tests/data/test_bsSeeker.Mouse.SRR892982_2.fastq tests/data/bsSeeker.Mouse.
↳SRR892982_2.fastq

```

(continues on next page)

(continued from previous page)

```

6
7 gzip tests/data/bsSeeker.Mouse.SRR892982_1.fastq
8 gzip tests/data/bsSeeker.Mouse.SRR892982_2.fastq

```

Test Scripts

The following are the tests for checking that the tools in the WGBS pipeline are functioning correctly.

The tests should be run in this order so that the required input files are generated at the correct stage.

```

1 pytest -m wgs tests/test_fastqc_validation.py
2 pytest -m wgs tests/test_bs_seeker_filter.py
3 pytest -m wgs tests/test_bs_seeker_indexer.py
4 pytest -m wgs tests/test_bs_seeker_aligner.py
5 pytest -m wgs tests/test_bs_seeker_methylation_caller.py

```

These can be called as part of a single tool chain with:

```

1 python tests/test_toolchains.py --pipeline wgs

```

10.1.6 Hi-C Test Data

Test Data

Dataset

Stable ID	SRR1658573
Citation	Rao et al 2014

Genome

Assembly	GRCh38
Chromosome	21
Start	15000000
End	19999999

Method

The full dataset was downloaded from ENA aligned to the genome using GEM.

```

1 wget ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR165/003/SRR1658573/SRR1658573_1.fastq.gz
2 wget ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR165/003/SRR1658573/SRR1658573_2.fastq.gz
3
4 wget "http://www.ebi.ac.uk/ena/data/view/CM000663.2,CM000664.2,CM000665.2,CM000666.2,
↳CM000667.2,CM000668.2,CM000669.2,CM000670.2,CM000671.2,CM000672.2,CM000673.2,
↳CM000674.2,CM000675.2,CM000676.2,CM000677.2,CM000678.2,CM000679.2,CM000680.2,
↳CM000681.2,CM000682.2,CM000683.2,CM000684.2,CM000685.2,CM000686.2,J01415.2&
↳display=fasta&download=fasta&filename=entry.fasta" -O GCA_000001405.22.fasta

```

```

1 from tool.common import common
2
3 genome_file = "GCA_000001405.22.fasta"
4 new_genome_file = "GCA_000001405.22_gem.fasta"
5
6 common_handle = common()
7 common_handle.replaceENAHeader(genome_file, new_genome_file)
8
9 idx_loc = common_handle.gem_index_genome(new_genome_file, new_genome_file)
10
11 from pytabit.mapping.mapper import full_mapping
12
13 fastq_file_1 = "SRR1658573_1.fastq"
14 fastq_file_2 = "SRR1658573_2.fastq"
15 output_dir = "mapped_reads"
16 windows = None
17 enzyme_name = "MboI"
18
19 map_files_1 = full_mapping(
20     idx_loc, fastq_file, output_dir,
21     r_enz=enzyme_name, windows=windows, frag_map=True, nthreads=32,
22     clean=True, temp_dir='/tmp/'
23 )
24 map_files_2 = full_mapping(
25     idx_loc, fastq_file, output_dir,
26     r_enz=enzyme_name, windows=windows, frag_map=True, nthreads=32,
27     clean=True, temp_dir='/tmp/'
28 )

```

A list of FASTQ ids that had aligned to the genome were then extracted:

```

1 grep chr21 mapped_reads/SRR1658573_1_full_1-100.map | tr "\t" "~" | cut -d"~" -f1 -f5
↪ | tr ":" "\t" | awk '(NR==1) || (($4>15000000) && ($4<20000000))' | tr "\t" "~" |
↪ cut -d "~" -f1 > SRR1658573_1_chr21_1-100.row
2 grep chr21 mapped_reads/SRR1658573_2_full_1-100.map | tr "\t" "~" | cut -d"~" -f1 -f5
↪ | tr ":" "\t" | awk '(NR==1) || (($4>15000000) && ($4<20000000))' | tr "\t" "~" |
↪ cut -d "~" -f1 > SRR1658573_2_chr21_1-100.row

```

The IDs were filtered to ensure matching pairs in both files:

```

1 grep -Fx -f SRR1658573_1_chr21_1-100.row SRR1658573_2_chr21_1-100.row > SRR1658573_
↪ chr21.row

```

The *split_paired_fastq.py* was used to divide the original FASTQ files into chunks of 1000000 reads. The *ExtractRowsFromFASTQ.py* script was then used to extract the matching FASTQ pairs from each of the FASTQ files in parallel. All of the individual FASTQ files were then concatenated together to form the final 2 FASTQ test files. The commands for this were:

```

1 python split_paired_fastq.py --input_1 SRR1658573_1.fastq --input_2 SRR1658573_1.fastq

```

Test Scripts

The following are the tests for checking that the tools in the Hi-C pipeline are functioning correctly.

The tests should be run in this order so that the required input files are generated at the correct stage.

```

1 pytest tests/test_gem_indexer.py
2 pytest tests/test_tb_full_mapping.py
3 pytest tests/test_tb_parse_mapping.py
4 pytest tests/test_tb_filter.py
5 pytest tests/test_tb_generate_tads.py
6 pytest tests/test_tb_save_hdf5_matrix.py

```

These can be called as part of a single tool chain with:

```

1 python tests/test_toolchains.py --pipeline hic

```

10.2 Pipelines

There is a test for each of the tools. This uses the “process” scripts to run each of the tools. This is to ensure that the pipeline scripts are able to call out to each of the tools and the correct parameters are handed to each one.

10.2.1 ChIP-Seq

To run the pipeline test:

```

pytest tests/test_pipeline_chipseq.py

```

Methods

`tests.test_pipeline_chipseq.test_chipseq_pipeline_00()`

Test case to ensure that the ChIP-seq pipeline code works.

Running the pipeline with the test data from the command line:

```

runcompss \
  --lang=python \
  --library_path=${HOME}/bin \
  --pythonpath=/<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
  --log_level=debug \
  process_chipseq.py \
  --taxon_id 9606 \
  --genome /<dataset_dir>/Human.GCA_000001405.22.fasta \
  --assembly GRCh38 \
  --file /<dataset_dir>/DRR000150.22.fastq

```

`tests.test_pipeline_chipseq.test_chipseq_pipeline_01()`

Test case to ensure that the ChIP-seq pipeline code works.

Running the pipeline with the test data from the command line:

```

runcompss \
  --lang=python \
  --library_path=${HOME}/bin \
  --pythonpath=/<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
  --log_level=debug \
  process_chipseq.py \
  --taxon_id 9606

```

(continues on next page)

(continued from previous page)

```
--genome /<dataset_dir>/Human.GCA_000001405.22.fasta      \  
--assembly GRCh38                                         \  
--file /<dataset_dir>/DRR000150.22.fastq
```

Sample Data

Test Data for ChIP-seq pipeline

10.2.2 iDamID-Seq

To run the pipeline test:

```
pytest tests/test_pipeline_idamidseq.py
```

Methods

`tests.test_pipeline_idamidseq.test_idamidseq_pipeline_00()`

Test case to ensure that the ChIP-seq pipeline code works.

Running the pipeline with the test data from the command line:

```
runcompss                                               \  
  --lang=python                                         \  
  --library_path=${HOME}/bin                           \  
  --pythonpath=/<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \  
  --log_level=debug                                     \  
process_damidseq.py                                     \  
  --taxon_id 9606                                       \  
  --genome /<dataset_dir>/Human.GCA_000001405.22.fasta \  
  --assembly GRCh38                                     \  
  --file /<dataset_dir>/DRR000150.22.fastq
```

`tests.test_pipeline_idamidseq.test_idamidseq_pipeline_01()`

Test case to ensure that the ChIP-seq pipeline code works.

Running the pipeline with the test data from the command line:

```
runcompss                                               \  
  --lang=python                                         \  
  --library_path=${HOME}/bin                           \  
  --pythonpath=/<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \  
  --log_level=debug                                     \  
process_damidseq.py                                     \  
  --taxon_id 9606                                       \  
  --genome /<dataset_dir>/Human.GCA_000001405.22.fasta \  
  --assembly GRCh38                                     \  
  --file /<dataset_dir>/DRR000150.22.fastq
```

Sample Data

Test Data for ChIP-seq pipeline

10.2.3 Genome Indexing

To run the pipeline test:

```
pytest tests/test_pipeline_genome.py
```

Methods

`tests.test_pipeline_genome.test_genome_pipeline_00()`

Test case to ensure that the Genome indexing pipeline code works.

Running the pipeline with the test data from the command line:

```
runcompss \
  --lang=python \
  --library_path=${HOME}/bin \
  --pythonpath=/<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
  --log_level=debug \
  process_genome.py \
    --taxon_id 9606 \
    --genome /<dataset_dir>/Human.GCA_000001405.22.fasta \
    --assembly GRCh38 \
    --file /<dataset_dir>/DRR000150.22.fastq
```

`tests.test_pipeline_genome.test_genome_pipeline_01()`

Test case to ensure that the Genome indexing pipeline code works.

Running the pipeline with the test data from the command line:

```
runcompss \
  --lang=python \
  --library_path=${HOME}/bin \
  --pythonpath=/<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
  --log_level=debug \
  process_genome.py \
    --taxon_id 9606 \
    --genome /<dataset_dir>/Human.GCA_000001405.22.fasta \
    --assembly GRCh38 \
    --file /<dataset_dir>/DRR000150.22.fastq
```

Sample Data

Uses the genome sequences required by all the tools and pipelines

10.2.4 Hi-C

To run the pipeline test:

```
pytest tests/test_pipeline_tb.py
```

Methods

tests.test_pipeline_tb.test_tb_pipeline()

Test case to ensure that the Hi-C pipeline code works.

Running the pipeline with the test data from the command line:

```
runcompss \
  --lang=python \
  --library_path=/home/compss/bin \
  --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
  --log_level=debug \
  process_hic.py \
  --taxon_id 9606 \
  --genome /<dataset_dir>/tb.Human.GCA_000001405.22_gem.fasta \
  --assembly GRCh38 \
  --file1 /<dataset_dir>/tb.Human.SRR1658573_1.fastq \
  --file2 /<dataset_dir>/tb.Human.SRR1658573_2.fastq \
  --genome_gem /<dataset_dir>/tb.Human.GCA_000001405.22_gem.fasta.gem \
  --taxon_id 9606 \
  --enzyme_name MboI \
  --resolutions 10000,100000 \
  --windows1 1,100 \
  --windows2 1,100 \
  --normalized 1 \
  --tag tb.Human.SRR1658573 \
  --window_type frag \
```

Sample Data

Hi-C Test Data

10.2.5 MNase-Seq

To run the pipeline test:

```
pytest tests/test_pipeline_mnaseseq.py
```

Methods

tests.test_pipeline_mnaseseq.test_mnaseseq_pipeline()

Test case to ensure that the MNase-seq pipeline code works.

Running the pipeline with the test data from the command line:

```
runcompss \
  --lang=python \
  --library_path=${HOME}/bin \
  --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
  --log_level=debug \
  process_mnaseseq.py \
  --taxon_id 10090 \
  --genome /<dataset_dir>/Mouse.GRCm38.fasta \
  --assembly GRCm38 \
  --file /<dataset_dir>/DRR000386.fastq \
```

Sample Data

Test Data for MNase-seq pipeline

10.2.6 RNA-Seq

To run the pipeline test:

```
pytest tests/test_pipeline_rnaseq.py
```

Methods

`tests.test_pipeline_rnaseq.test_rnaseq_pipeline()`

Test case to ensure that the RNA-seq pipeline code works.

Running the pipeline with the test data from the command line:

```
runcompss \
  --lang=python \
  --library_path=${HOME}/bin \
  --pythonpath=/<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \
  --log_level=debug \
  process_rnaseq.py \
  --taxon_id 9606 \
  --genome /<dataset_dir>/Human.GRCh38.fasta \
  --assembly GRCh38 \
  --file /<dataset_dir>/ERR030872_1.fastq \
  --file2 /<dataset_dir>/ERR030872_2.fastq
```

Sample Data

Test Data for RNA-seq pipeline

10.2.7 Whole Genome Bisulfate Sequencing (WGBS)

To run the pipeline test:

```
pytest tests/test_pipeline_wgbs.py
```

Methods

Sample Data

WGBS Test Data

10.3 Tools

As the data stored is only the raw data, each of the sets of tools has been packaged up into a tool chain to run each of the tools without failing. This has been done with the `tests.test_toolchains.py` script.

```
python tests/test_toolchains.py --pipeline [genome | chipseq | hic | mnaseq |   
↪rnaseq | wgs]
```

This script automates the running of each of the tools that are required for a given pipeline.

10.3.1 Methods

`tests.test_toolchains.all_toolchain(verbose=False)`

Runs the tests for all of the tools

This set is only required for determining code coverage.

`tests.test_toolchains.biobambam_toolchain(verbose=False)`

Runs the tests for all of the tools from the BWA pipeline

Runs the following tests:

```
pytest -m chipseq tests/test_fastqc_validation.py  
pytest -m chipseq tests/test_bwa_indexer.py  
pytest -m chipseq tests/test_bwa_aligner.py  
pytest -m chipseq tests/test_biobambam.py
```

`tests.test_toolchains.bowtie2_toolchain(verbose=False)`

Runs the tests for all of the tools from the BWA pipeline

Runs the following tests:

```
pytest -m bowtie2 tests/test_fastqc_validation.py  
pytest -m bowtie2 tests/test_bowtie_indexer.py  
pytest -m bowtie2 tests/test_bowtie2_aligner.py
```

`tests.test_toolchains.bwa_toolchain(verbose=False)`

Runs the tests for all of the tools from the BWA pipeline

Runs the following tests:

```
pytest -m bwa tests/test_fastqc_validation.py  
pytest -m bwa tests/test_bwa_indexer.py  
pytest -m bwa tests/test_bwa_aligner.py
```

`tests.test_toolchains.chipseq_toolchain(verbose=False)`

Runs the tests for all of the tools from the ChIP-seq pipeline

Runs the following tests:

```
pytest -m chipseq tests/test_fastqc_validation.py  
pytest -m chipseq tests/test_bwa_indexer.py  
pytest -m chipseq tests/test_bwa_aligner.py  
pytest -m chipseq tests/test_biobambam.py  
pytest -m chipseq tests/test_mac2.py
```

`tests.test_toolchains.genome_toolchain(verbose=False)`

Runs the tests for all of the tools from the Genome indexing pipeline

Runs the following tests:


```

pytest -m genome tests/test_bowtie_indexer.py
pytest -m genome tests/test_bwa_indexer.py
pytest -m genome tests/test_gem_indexer.py

```

`tests.test_toolchains.hic_toolchain` (*verbose=False*)

Runs the tests for all of the tools from the Hi-C pipeline

Runs the following tests:

```

pytest -m hic tests/test_fastqc_validation.py
pytest -m hic tests/test_gem_indexer.py
pytest -m hic tests/test_tb_full_mapping.py
pytest -m hic tests/test_tb_parse_mapping.py
pytest -m hic tests/test_tb_filter.py
pytest -m hic tests/test_tb_normalize.py
pytest -m hic tests/test_tb_segment.py
pytest -m hic tests/test_tb_generate_tads.py
pytest -m hic tests/test_tb_bin.py
pytest -m hic tests/test_tb_save_hdf5_matrix.py

```

`tests.test_toolchains.idamidseq_toolchain` (*verbose=False*)

Runs the tests for all of the tools from the iDamID-seq pipeline

Runs the following tests:

```

pytest -m idamidseq tests/test_bwa_indexer.py
pytest -m idamidseq tests/test_bwa_aligner.py
pytest -m idamidseq tests/test_biobambam.py
pytest -m idamidseq tests/test_hsgenome.py
pytest -m idamidseq tests/test_idear.py

```

`tests.test_toolchains.mnaseq_toolchain` (*verbose=False*)

Runs the tests for all of the tools from the MNase-seq pipeline

Runs the following tests:

```

pytest -m mnaseq tests/test_fastqc_validation.py
pytest -m mnaseq tests/test_bwa_indexer.py
pytest -m mnaseq tests/test_bwa_aligner.py
pytest -m mnaseq tests/test_inps.py

```

`tests.test_toolchains.rnaseq_toolchain` (*verbose=False*)

Runs the tests for all of the tools from the RNA-seq pipeline

Runs the following tests:

```

pytest -m rnaseq tests/test_fastqc_validation.py
pytest -m rnaseq tests/test_kallisto_indexer.py
pytest -m rnaseq tests/test_kallisto_quant.py

```

`tests.test_toolchains.tidy_data` ()

Runs the tidy_data.sh script

`tests.test_toolchains.wgbs_toolchain` (*verbose=0*)

Runs the tests for all of the tools from the WGBS pipeline

Runs the following tests:

```
pytest -m wgs tests/test_fastqc_validation.py
pytest -m wgs tests/test_bs_seeker_filter.py
pytest -m wgs tests/test_bs_seeker_indexer.py
pytest -m wgs tests/test_bs_seeker_aligner.py
pytest -m wgs tests/test_bs_seeker_methylation_caller.py
```

CHAPTER 11

Indices and tables

- `genindex`
- `modindex`
- `search`

p

process_align_bowtie, 16
process_align_bwa, 23
process_align_bwa_mem, 24
process_biobambam, 14
process_bs_seeker_aligner, 20
process_bs_seeker_filter, 21
process_bs_seeker_index, 18
process_bsgenome, 17
process_chipseq, 26
process_damidseq, 28
process_genome, 13
process_hic, 36
process_iNPS, 29
process_mac2, 30
process_mnaseq, 32
process_rnaseq, 33
process_trim_galore, 35

t

tests, 109
tests.test_pipeline_chipseq, 109
tests.test_pipeline_genome, 111
tests.test_pipeline_idamidseq, 110
tests.test_pipeline_mnaseq, 112
tests.test_pipeline_rnaseq, 113
tests.test_pipeline_tb, 112
tests.test_toolchains, 114
tool, 39

A

alignerUtils (class in *tool.aligner_utils*), 71
 all_toolchain() (in module *tests.test_toolchains*),
 114

B

bam_copy() (*tool.bam_utils.bamUtils* static method),
 73
 bam_copy() (*tool.bam_utils.bamUtilsTask* method), 75
 bam_count_reads() (*tool.bam_utils.bamUtils* static
 method), 73
 bam_filter() (*tool.bam_utils.bamUtils* static
 method), 73
 bam_filter() (*tool.bam_utils.bamUtilsTask* method),
 75
 bam_index() (*tool.bam_utils.bamUtils* static method),
 74
 bam_index() (*tool.bam_utils.bamUtilsTask* method),
 75
 bam_list_chromosomes()
 (*tool.bam_utils.bamUtils* static method),
 74
 bam_list_chromosomes()
 (*tool.bam_utils.bamUtilsTask* method), 75
 bam_merge() (*tool.bam_utils.bamUtils* static method),
 74
 bam_merge() (*tool.bam_utils.bamUtilsTask* method),
 75
 bam_merge_10() (*tool.bam_utils.bamUtilsTask*
 method), 75
 bam_merge_2() (*tool.bam_utils.bamUtilsTask*
 method), 76
 bam_merge_3() (*tool.bam_utils.bamUtilsTask*
 method), 76
 bam_merge_4() (*tool.bam_utils.bamUtilsTask*
 method), 76
 bam_merge_5() (*tool.bam_utils.bamUtilsTask*
 method), 76
 bam_paired_reads() (*tool.bam_utils.bamUtils*
 static method), 74
 bam_paired_reads() (*tool.bam_utils.bamUtilsTask*
 method), 76
 bam_sort() (*tool.bam_utils.bamUtils* static method),
 74
 bam_sort() (*tool.bam_utils.bamUtilsTask* method), 77
 bam_split() (*tool.bam_utils.bamUtils* static method),
 74
 bam_stats() (*tool.bam_utils.bamUtils* static method),
 74
 bam_stats() (*tool.bam_utils.bamUtilsTask* method),
 77
 bam_to_bed() (*tool.bam_utils.bamUtils* static
 method), 74
 bamUtils (class in *tool.bam_utils*), 73
 bamUtilsTask (class in *tool.bam_utils*), 75
 biobambam (class in *tool.biobambam_filter*), 51
 biobambam_filter_alignments()
 (*tool.biobambam_filter.biobambam* method),
 51
 biobambam_toolchain() (in module
tests.test_toolchains), 114
 bowtie2_align_reads()
 (*tool.aligner_utils.alignerUtils* static method),
 71
 bowtie2_aligner_paired()
 (*tool.bowtie_aligner.bowtie2AlignerTool*
 method), 45
 bowtie2_aligner_single()
 (*tool.bowtie_aligner.bowtie2AlignerTool*
 method), 45
 bowtie2_indexer()
 (*tool.bowtie_indexer.bowtieIndexerTool*
 method), 41
 bowtie2_toolchain() (in module
tests.test_toolchains), 114
 bowtie2_untar_index()
 (*tool.aligner_utils.alignerUtils* method),
 71
 bowtie2AlignerTool (class in *tool.bowtie_aligner*),

45
 bowtie_index_genome() (tool.aligner_utils.alignerUtils static method), 72
 bowtieIndexerTool (class in tool.bowtie_indexer), 41
 bs_seeker_aligner() (tool.bs_seeker_aligner.bssAlignerTool method), 50
 bs_seeker_aligner_single() (tool.bs_seeker_aligner.bssAlignerTool method), 50
 bsgenome_creator() (tool.forge_bsgenome.bsgenomeTool method), 41
 bsgenomeTool (class in tool.forge_bsgenome), 41
 bss_build_index() (tool.bs_seeker_indexer.bssIndexerTool method), 42
 bss_seeker_filter() (tool.bs_seeker_filter.filterReadsTool method), 52
 bssAlignerTool (class in tool.bs_seeker_aligner), 50
 bssIndexerTool (class in tool.bs_seeker_indexer), 42
 bwa_aligner_paired() (tool.bwa_aligner.bwaAlignerTool method), 47
 bwa_aligner_paired() (tool.bwa_mem_aligner.bwaAlignerMEMTool method), 48
 bwa_aligner_single() (tool.bwa_aligner.bwaAlignerTool method), 47
 bwa_aligner_single() (tool.bwa_mem_aligner.bwaAlignerMEMTool method), 49
 bwa_aln_align_reads_paired() (tool.aligner_utils.alignerUtils method), 72
 bwa_aln_align_reads_single() (tool.aligner_utils.alignerUtils method), 72
 bwa_index_genome() (tool.aligner_utils.alignerUtils static method), 72
 bwa_indexer() (tool.bwa_indexer.bwaIndexerTool method), 43
 bwa_mem_align_reads() (tool.aligner_utils.alignerUtils static method), 72
 bwa_toolchain() (in module tests.test_toolchains), 114
 bwa_untar_index() (tool.aligner_utils.alignerUtils method), 73
 bwaAlignerMEMTool (class in tool.bwa_mem_aligner), 48
 bwaAlignerTool (class in tool.bwa_aligner), 47
 bwaIndexerTool (class in tool.bwa_indexer), 43

C

cd (class in tool.common), 71
 check_header() (tool.bam_utils.bamUtils static method), 75
 check_header() (tool.bam_utils.bamUtilsTask method), 77
 chipseq_toolchain() (in module tests.test_toolchains), 114
 closeFastQ() (tool.fastqreader.fastqreader method), 77
 closeOutputFiles() (tool.fastqreader.fastqreader method), 77
 createOutputFiles() (tool.fastqreader.fastqreader method), 77

E

eof() (tool.fastqreader.fastqreader method), 77

F

fastq_match_paired_ends() (tool.fastq_utils.fastqUtils static method), 79
 fastq_randomise() (tool.fastq_utils.fastqUtils static method), 79
 fastq_sort_file() (tool.fastq_utils.fastqUtils static method), 79
 fastq_splitter (class in tool.fastq_splitter), 78
 fastqcTool (class in tool.validate_fastqc), 39
 fastqreader (class in tool.fastqreader), 77
 fastqUtils (class in tool.fastq_utils), 79
 filterReadsTool (class in tool.bs_seeker_filter), 52

G

gem_index_genome() (tool.aligner_utils.alignerUtils static method), 73
 gem_indexer() (tool.gem_indexer.gemIndexerTool method), 44
 gemIndexerTool (class in tool.gem_indexer), 44
 genome_to_2bit() (tool.forge_bsgenome.bsgenomeTool static method), 42
 genome_toolchain() (in module tests.test_toolchains), 114
 get_aln_params() (tool.bowtie_aligner.bowtie2AlignerTool static method), 46
 get_aln_params() (tool.bs_seeker_aligner.bssAlignerTool static method), 51
 get_aln_params() (tool.bwa_aligner.bwaAlignerTool static method), 47

[get_bss_index_params\(\)](#) (*tool.bs_seeker_indexer.bssIndexerTool* static method), 43
[get_chrom_size\(\)](#) (*tool.forge_bsgenome.bsgenomeTool* static method), 42
[get_mac2_params\(\)](#) (*tool.mac2.mac2* static method), 56
[get_mem_params\(\)](#) (*tool.bwa_mem_aligner.bwaAlignerMEMTool* static method), 49
[get_trimgalore_params\(\)](#) (*tool.trimgalore.trimgalore* static method), 40, 53

H

[hic_toolchain\(\)](#) (in module *tests.test_toolchains*), 115

I

[idamidseq_toolchain\(\)](#) (in module *tests.test_toolchains*), 115
[idear_peak_calling\(\)](#) (*tool.idear.idearTool* method), 54
idearTool (class in *tool.idear*), 54
[incrementOutputFiles\(\)](#) (*tool.fastqreader.fastqreader* method), 77
inps (class in *tool.inps*), 55
[inps_peak_calling\(\)](#) (*tool.inps.inps* method), 55

K

[kallisto_indexer\(\)](#) (*tool.kallisto_indexer.kallistoIndexerTool* method), 44
[kallisto_quant_paired\(\)](#) (*tool.kallisto_quant.kallistoQuantificationTool* method), 55
[kallisto_quant_single\(\)](#) (*tool.kallisto_quant.kallistoQuantificationTool* method), 55
[kallisto_tsv2bed\(\)](#) (*tool.kallisto_quant.kallistoQuantificationTool* method), 56
[kallisto_tsv2gff\(\)](#) (*tool.kallisto_quant.kallistoQuantificationTool* method), 56
kallistoIndexerTool (class in *tool.kallisto_indexer*), 44
kallistoQuantificationTool (class in *tool.kallisto_quant*), 55

L

[load_gff_ensembl\(\)](#) (*tool.kallisto_quant.kallistoQuantificationTool* static method), 56
[load_gff_ucsc\(\)](#) (*tool.kallisto_quant.kallistoQuantificationTool* static method), 56

M

mac2 (class in *tool.mac2*), 56
[mac2_peak_calling\(\)](#) (*tool.mac2.mac2* method), 56
[mac2_peak_calling_nobgd\(\)](#) (*tool.mac2.mac2* method), 57
[mnaseseq_toolchain\(\)](#) (in module *tests.test_toolchains*), 115

N

[next\(\)](#) (*tool.fastqreader.fastqreader* method), 77

O

[openFastQ\(\)](#) (*tool.fastqreader.fastqreader* method), 78

P

[paired_splitter\(\)](#) (*tool.fastq_splitter.fastq_splitter* method), 78
process_align_bowtie (module), 16
process_align_bwa (module), 23
process_align_bwa_mem (module), 24
process_biobambam (class in *process_biobambam*), 15
process_biobambam (module), 14
process_bowtie (class in *process_align_bowtie*), 17
process_bs_seeker_aligner (class in *process_bs_seeker_aligner*), 21
process_bs_seeker_aligner (module), 20
process_bs_seeker_filter (module), 21
process_bs_seeker_index (class in *process_bs_seeker_index*), 19
process_bs_seeker_index (module), 18
process_bsFilter (class in *process_bs_seeker_filter*), 22
process_bsgenome (class in *process_bsgenome*), 18
process_bsgenome (module), 17
process_bwa (class in *process_align_bwa*), 24
process_bwa_mem (class in *process_align_bwa_mem*), 25
process_chipseq (class in *process_chipseq*), 27
process_chipseq (module), 26
process_damidseq (class in *process_damidseq*), 28
process_damidseq (module), 28
process_genome (class in *process_genome*), 14
process_genome (module), 13
process_hic (class in *process_hic*), 38
process_hic (module), 36
process_iNPS (class in *process_iNPS*), 30
process_iNPS (module), 29

process_mac2 (class in process_mac2), 31
 process_mac2 (module), 30
 process_mnaseq (class in process_mnaseq), 33
 process_mnaseq (module), 32
 process_rnaseq (class in process_rnaseq), 34
 process_rnaseq (module), 33
 process_trim_galore (class in process_trim_galore), 36
 process_trim_galore (module), 35

R

replaceENAHeader ()
 (tool.aligner_utils.alignerUtils static method), 73
 rnaseq_toolchain () (in module tests.test_toolchains), 115
 run () (process_align_bowtie.process_bowtie method), 17
 run () (process_align_bwa.process_bwa method), 24
 run () (process_align_bwa_mem.process_bwa_mem method), 25
 run () (process_biobambam.process_biobambam method), 15
 run () (process_bs_seeker_aligner.process_bs_seeker_aligner method), 21
 run () (process_bs_seeker_filter.process_bsFilter method), 22
 run () (process_bs_seeker_index.process_bs_seeker_index method), 19
 run () (process_bsgenome.process_bsgenome method), 18
 run () (process_chipseq.process_chipseq method), 27
 run () (process_damidseq.process_damidseq method), 28
 run () (process_genome.process_genome method), 14
 run () (process_hic.process_hic method), 38
 run () (process_iNPS.process_iNPS method), 30
 run () (process_mac2.process_mac2 method), 31
 run () (process_mnaseq.process_mnaseq method), 33
 run () (process_rnaseq.process_rnaseq method), 34
 run () (process_trim_galore.process_trim_galore method), 36
 run () (tool.biobambam_filter.biobambam method), 52
 run () (tool.bowtie_aligner.bowtie2AlignerTool method), 46
 run () (tool.bowtie_indexer.bowtieIndexerTool method), 41
 run () (tool.bs_seeker_aligner.bssAlignerTool method), 51
 run () (tool.bs_seeker_filter.filterReadsTool method), 52
 run () (tool.bs_seeker_indexer.bssIndexerTool method), 43
 run () (tool.bwa_aligner.bwaAlignerTool method), 48

run () (tool.bwa_indexer.bwaIndexerTool method), 43
 run () (tool.bwa_mem_aligner.bwaAlignerMEMTool method), 49
 run () (tool.fastq_splitter.fastq_splitter method), 78
 run () (tool.forge_bsgenome.bsgenomeTool method), 42
 run () (tool.gem_indexer.gemIndexerTool method), 44
 run () (tool.idear.idearTool method), 54
 run () (tool.inps.inps method), 55
 run () (tool.kallisto_indexer.kallistoIndexerTool method), 44
 run () (tool.kallisto_quant.kallistoQuantificationTool method), 56
 run () (tool.mac2.mac2 method), 58
 run () (tool.tb_bin.tbBinTool method), 64
 run () (tool.tb_filter.tbFilterTool method), 62
 run () (tool.tb_full_mapping.tbFullMappingTool method), 58
 run () (tool.tb_generate_tads.tbGenerateTADsTool method), 66
 run () (tool.tb_model.tbModelTool method), 67
 run () (tool.tb_normalize.tbNormalizeTool method), 63
 run () (tool.tb_parse_mapping.tbParseMappingTool method), 59
 run () (tool.tb_save_hdf5_matrix.tbSaveAdjacencyHDF5Tool method), 66
 run () (tool.tb_segment.tbSegmentTool method), 62
 run () (tool.trim_galore.trim_galore method), 40, 53
 run () (tool.validate_fastqc.fastqcTool method), 39
 run_aligner () (tool.bs_seeker_aligner.bssAlignerTool method), 51

S

sam_to_bam () (tool.bam_utils.bamUtils static method), 75
 seq_read_stats () (tool.kallisto_quant.kallistoQuantificationTool static method), 56
 single_splitter ()
 (tool.fastq_splitter.fastq_splitter method), 79

T

tb_bin () (tool.tb_bin.tbBinTool method), 65
 tb_filter () (tool.tb_filter.tbFilterTool method), 62
 tb_full_mapping_frag ()
 (tool.tb_full_mapping.tbFullMappingTool method), 58
 tb_full_mapping_iter ()
 (tool.tb_full_mapping.tbFullMappingTool method), 59
 tb_generate_tads ()
 (tool.tb_generate_tads.tbGenerateTADsTool method), 67
 tb_hic_chr () (tool.tb_generate_tads.tbGenerateTADsTool method), 67

tb_matrix_hdf5() (*tool.tb_save_hdf5_matrix.tbSaveAdjacencyHDF5Tool* method), 66
tb_merge_tad_files() (*tool.tb_generate_tads.tbGenerateTADsTool* method), 67
tb_model() (*tool.tb_model.tbModelTool* method), 68
tb_normalize() (*tool.tb_normalize.tbNormalizeTool* method), 64
tb_parse_mapping_frag() (*tool.tb_parse_mapping.tbParseMappingTool* method), 61
tb_parse_mapping_iter() (*tool.tb_parse_mapping.tbParseMappingTool* method), 61
tb_segment() (*tool.tb_segment.tbSegmentTool* method), 63
tbBinTool (class in *tool.tb_bin*), 64
tbFilterTool (class in *tool.tb_filter*), 62
tbFullMappingTool (class in *tool.tb_full_mapping*), 58
tbGenerateTADsTool (class in *tool.tb_generate_tads*), 66
tbModelTool (class in *tool.tb_model*), 67
tbNormalizeTool (class in *tool.tb_normalize*), 63
tbParseMappingTool (class in *tool.tb_parse_mapping*), 59
tbSaveAdjacencyHDF5Tool (class in *tool.tb_save_hdf5_matrix*), 66
tbSegmentTool (class in *tool.tb_segment*), 62
test_chipseq_pipeline_00() (in module *tests.test_pipeline_chipseq*), 109
test_chipseq_pipeline_01() (in module *tests.test_pipeline_chipseq*), 109
test_genome_pipeline_00() (in module *tests.test_pipeline_genome*), 111
test_genome_pipeline_01() (in module *tests.test_pipeline_genome*), 111
test_idamidseq_pipeline_00() (in module *tests.test_pipeline_idamidseq*), 110
test_idamidseq_pipeline_01() (in module *tests.test_pipeline_idamidseq*), 110
test_mnaseseq_pipeline() (in module *tests.test_pipeline_mnaseseq*), 112
test_rnaseq_pipeline() (in module *tests.test_pipeline_rnaseq*), 113
test_tb_pipeline() (in module *tests.test_pipeline_tb*), 112
tests (module), 109
tests.test_pipeline_chipseq (module), 109
tests.test_pipeline_genome (module), 111
tests.test_pipeline_idamidseq (module), 110
tests.test_pipeline_mnaseseq (module), 112
tests.test_pipeline_rnaseq (module), 113
tests.test_toolchains (module), 114
tidy_data() (in module *tests.test_toolchains*), 115
tool (module), 39, 71
trimgalore (class in *tool.trimgalore*), 40, 53
trimgalore_paired() (*tool.trimgalore.trimgalore* method), 40, 53
trimgalore_single() (*tool.trimgalore.trimgalore* method), 40, 53
trimgalore_version() (*tool.trimgalore.trimgalore* method), 41, 53
U
untar_index() (*tool.bowtie_aligner.bowtie2AlignerTool* method), 46
untar_index() (*tool.bwa_aligner.bwaAlignerTool* method), 48
untar_index() (*tool.bwa_mem_aligner.bwaAlignerMEMTool* method), 49
V
validate() (*tool.validate_fastqc.fastqcTool* method), 39
W
wgbs_toolchain() (in module *tests.test_toolchains*), 115
writeOutput() (*tool.fastqreader.fastqreader* method), 78