

---

# **MuG - CHi-C Pipeline Documentation**

***Release 0.1***

**Pablo Acera**

**Feb 25, 2019**



---

## Table of Contents

---

<b>1</b>	<b>Requirements and Installation</b>	<b>1</b>
1.1	Requirements . . . . .	1
1.2	Installation . . . . .	2
<b>2</b>	<b>Full Installation</b>	<b>3</b>
2.1	Setup the System Environment . . . . .	3
2.2	Setup pyenv and pyenv-virtualenv . . . . .	3
2.3	Installation Process . . . . .	4
2.4	Setup the symlinks . . . . .	5
2.5	Prepare the Python Environment . . . . .	6
<b>3</b>	<b>Pipelines</b>	<b>7</b>
3.1	Map and parse CHi-C reads . . . . .	7
3.2	Create CHiCAGO input RMAP . . . . .	8
3.3	Create CHiCAGO input BAITMAP . . . . .	8
3.4	Create CHiCAGO input Design files . . . . .	8
3.5	Convert BAM file into chicago input files .chinput . . . . .	10
3.6	Data normalization and peak calling . . . . .	11
3.7	Run the entire CHi-C pipeline . . . . .	12
<b>4</b>	<b>Tools for processing fastq C-HiC files</b>	<b>13</b>
4.1	Map and parser reads . . . . .	13
4.2	Create CHiCAGO input files . . . . .	15
4.3	Convert bam files into chicago input . . . . .	16
4.4	Normalize data and call C-HiC peaks . . . . .	17
<b>5</b>	<b>Architectural Design Record</b>	<b>19</b>
5.1	25-09-2018 handling_chr_header branch merge with master . . . . .	19
5.2	15-10-2018 mm_mods_for_makebaitmaps branch merge with master . . . . .	19
5.3	10-12-2018 creation of the branch VM_CR1 (VM current release version) . . . . .	19
<b>6</b>	<b>License</b>	<b>21</b>
<b>7</b>	<b>Indices and tables</b>	<b>25</b>
	<b>Python Module Index</b>	<b>27</b>



---

## Requirements and Installation

---

### 1.1 Requirements

#### 1.1.1 Software

- Python 2.7.12+
- R >=3.1.2
- bedtools
- perl
- HiCUP
- bwa
- GEM
- TADbit
- samtools>1.3

#### 1.1.2 Python Modules

- mg-tool-api
- pylint
- pytest
- pandas
- rtree

### 1.1.3 R Modules

- argparse
- devtools
- Chicago

To Run `runChicago.py` and `process_runChicago.py`, the R script `runChicago.R` from <https://bitbucket.org/chicagoTeam/chicago/src/ceffddda8ea392a1e84e4db9593f8fc35ac88048/chicagoTools/?at=master> should be downloaded and added to PATH.

## 1.2 Installation

Directly from GitHub:

```
git clone https://github.com/pabloacera/C-HiC.git
```

Using pip:

```
pip install git+https://github.com/pabloacera/C-HiC.git
```

Install R modules, use the following R code:

```
install.packages("argparser")      install.packages("devtools")      library(devtools)      in-  
stall_bitbucket("chicagoTeam/Chicago", subdir="Chicago")
```

# CHAPTER 2

## Full Installation

The following document is for the full installation of all software required by the C-HiC module and all programmes that it uses. The document has been written with Ubuntu Linux, 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  
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 C-HiC
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 bedtools and libspatialindex-dev

```
1 sudo apt-get install bedtools
2 sudo apt-get install libspatialindex-dev
```

### 2.3.2 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
```

### 2.3.3 HiCUP

```
1 cd ${HOME}/lib
2 wget https://www.bioinformatics.babraham.ac.uk/projects/hicup/hicup_v0.6.1.tar.gz
3 tar -xzf hicup_v0.6.1.tar.gz
4 cd hicup_v0.6.1
5 chmod a+x *
```

### 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 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

```

### 2.3.6 Install CHiCAGO

```

1 sudo apt-get update -qq
2 sudo apt-get install python-rpy2
3
4
5 cd ${HOME}/lib
6 sudo apt-get install libtbb-dev
7 cd ${HOME}/C-HiC/
8 echo "R_LIB=${HOME}/R" > ${HOME}/.Renviro
9 echo ".libPaths('~ /R')" >> ${HOME}/.Rprofile
10 echo 'message("Using library:", .libPaths()[1])' >> ${HOME}/.Rprofile
11 R
12 options(repos = c(CRAN = "http://cran.rstudio.com"))
13 install.packages("Delaporte")
14 install.packages("MASS")
15
16 cd ${HOME}/C-HiC/CHiC/tool/scripts/
17 wget https://bitbucket.org/chicagoTeam/chicago/raw/
   ↪ e288015f75d36c5367d1595e0ac8099f2ce82aal/chicagoTools/bam2chicago.sh
18 chmod +x bam2chicago.sh

```

## 2.4 Setup the symlinks

```

1 cd ${HOME}/bin
2
3
4
5 ln -s ${HOME}/lib/hicup_v0.6.1/* ${HOME}/bin/
6
7 ln -s ${HOME}/lib/bwa/bwa bwa
8
9 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2 bowtie2

```

(continues on next page)

(continued from previous page)

```

10 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-align-s bowtie2-align-s
11 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-align-l bowtie2-align-l
12 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-build bowtie2-build
13 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-build-s bowtie2-build-s
14 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-build-l bowtie2-build-l
15 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-inspect bowtie2-inspect
16 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-inspect-s bowtie2-inspect-s
17 ln -s ${HOME}/lib/bowtie2-2.3.4-linux-x86_64/bowtie2-inspect-l bowtie2-inspect-l
18
19 ln -s ${HOME}/lib/htslib/bin/bgzip bgzip
20 ln -s ${HOME}/lib/htslib/bin/htsfile htsfile
21 ln -s ${HOME}/lib/htslib/bin/tabix tabix
22
23
24 ln -s ${HOME}/lib/samtools/bin/ace2sam ace2sam
25 ln -s ${HOME}/lib/samtools/bin/blast2sam.pl blast2sam.pl
26 ln -s ${HOME}/lib/samtools/bin/bowtie2sam.pl bowtie2sam.pl
27 ln -s ${HOME}/lib/samtools/bin/export2sam.pl export2sam.pl
28 ln -s ${HOME}/lib/samtools/bin/interpolate_sam.pl interpolate_sam.pl
29 ln -s ${HOME}/lib/samtools/bin/maq2sam-long maq2sam-long
30 ln -s ${HOME}/lib/samtools/bin/maq2sam-short maq2sam-short
31 ln -s ${HOME}/lib/samtools/bin/md5fa md5fa
32 ln -s ${HOME}/lib/samtools/bin/md5sum-lite md5sum-lite
33 ln -s ${HOME}/lib/samtools/bin/novo2sam.pl novo2sam.pl
34 ln -s ${HOME}/lib/samtools/bin/plot-bamstats plot-bamstats
35 ln -s ${HOME}/lib/samtools/bin/psl2sam.pl psl2sam.pl
36 ln -s ${HOME}/lib/samtools/bin/sam2vcf.pl sam2vcf.pl
37 ln -s ${HOME}/lib/samtools/bin/samtools samtools
38 ln -s ${HOME}/lib/samtools/bin/samtools.pl samtools.pl
39 ln -s ${HOME}/lib/samtools/bin/seq_cache_populate.pl seq_cache_populate.pl
40 ln -s ${HOME}/lib/samtools/bin/soap2sam.pl soap2sam.pl
41 ln -s ${HOME}/lib/samtools/bin/varfilter.py varfilter.py
42 ln -s ${HOME}/lib/samtools/bin/wgsim wgsim
43 ln -s ${HOME}/lib/samtools/bin/wgsim_eval.pl wgsim_eval.pl
44 ln -s ${HOME}/lib/samtools/bin/zoom2sam.pl zoom2sam.pl

```

## 2.5 Prepare the Python Environment

### 2.5.1 Install APIs and Pipelines

Checkout the code for the DM API and the C-HiC pipelines:

```

1 cd ${HOME}/code
2 pyenv activate C-HiC
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 pip install git+https://github.com/Multiscale-Genomics/mg-process-fastq.git
6
7
8 git clone https://github.com/pabloacera/C-HiC.git
9 cd C-HiC
10 pip install -e .
11 pip install -r requirements.txt
12 pip install dill

```

## 3.1 Map and parse CHi-C reads

This pipeline will take as input two fastq files, RE sites, the genome indexed with GEM and the same genome in FASTA file. This pipeline uses TADbit to map, filter and produce a bed file that will be used later on to produce bam file compatible with CHiCAGO algorithm. More information about filtering and mapping <https://3dgenomes.github.io/TADbit/>

### 3.1.1 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

**Wd** [folders and files] path to the working directory where the output files are

#### Example

REQUIREMENT - Needs two fastq files single end, FASTA genome and bowtie2 indexed genome.

When running the pipeline on a local machine without COMPSs:

```
1 python process_hicup.py \  
2   --config tests/json/config_hicup.json \  
3   --in_metadata tests/json/input_hicup.json \  
4
```

(continues on next page)

(continued from previous page)

```

4  --out_metadata tests/json/output_hicup.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_fastq2bed.py \
7  --config tests/json/config_hicup.json \
8  --in_metadata tests/json/input_hicup.json \
9  --out_metadata tests/json/output_hicup.json

```

### 3.1.2 Methods

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

This class run hicup tool which run hicup, doing the alignment and filtering of the reads and convert them into a BAM file.

**run** (*input\_files, metadata, output\_files*)

This is the main function that runs

#### Parameters

- **input\_files** (*dict*) – fastq1 fastq2
- **metadata** (*dict*) –
- **output\_files** (*dict*) –
- **out\_dir**: **str** directory to write the output

#### Returns

- **results** (*bool*)
- **output\_metadata** (*dict*)

## 3.2 Create CHiCAGO input RMAP

## 3.3 Create CHiCAGO input BAITMAP

## 3.4 Create CHiCAGO input Design files

This script use as input .rmap and .baitmap files and generate the Design files. NPerBin file (.npb): <baitID> <Total no. valid restriction fragments in distance bin 1> ... <Total no. valid restriction fragments in distance bin N>, where the bins map within the “proximal” distance range from each bait (0; maxLBrownEst] and bin size is defined by the binsize parameter. NBaitsPerBin file (.nbpb): <otherEndID> <Total no. valid baits in distance bin 1> ... <Total no. valid baits in distance bin N>, where the bins map within the “proximal” distance range from each other end (0; maxLBrownEst] and bin size is defined by the binsize parameter. Proximal Other End (ProxOE) file (.poe): <baitID> <otherEndID> <absolute distance> for all combinations of baits and other ends that map within the “proximal” distance range from

each other (0; maxLBrownEst]. Data in each file is preceded by a comment line listing the input parameters used to generate them.

### 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

“nbpb” : .nbpb file “npb” : .npb file “poe” : .poe file

#### Example

REQUIREMENT - Needs RMAP and BAITMAP files

When running the pipeline on a local machine without COMPSs:

```
1 python process_design.py \
2   --config tests/json/config_design.json \
3   --in_metadata tests/json/input_design.json \
4   --out_metadata tests/json/output_design.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_design.py   \
7   --config tests/json/config_design.json \
8   --in_metadata tests/json/input_design.json \
9   --out_metadata tests/json/output_design.json
```

### 3.4.2 Methods

**class** process\_design.**process\_design** (*configuration=None*)

This class generates the Design files and chinput files, input for CHiCAGO. Starting from rmap and baitmap and capture HiC BAM files.

**run** (*input\_files, metadata, output\_files*)

Main function to run the tools, MakeDesignFiles\_Tool.py and bam2chicago\_Tool.py

#### Parameters

- **input\_files** (*dict*) – designDir: path to the folder with .rmap and .baitmap files  
rmapFile: path to the .rmap file baitmapFile: path to the .baitmap file bamFile: path to the capture HiC bamfiles
- **metadata** (*dict*) – input metadata
- **output\_files** (*dict*) – outPrefixDesign : Path and name of the output prefix, recommend to be the same as rmap and baitmap files. sample\_name: Path and name of the .chinput file

#### Returns

- *bool*
- *output\_metadata*

## 3.5 Convert BAM file into chicago input files .chinput

This pipeline convert the output of process\_bed2bam.py BAM file to a .chinput file, input for process\_runChicago.py

### 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

chrRMAP : .rmap file with chr# format chrBAITMAP : .baitmap file with chr# format sample\_name : .chinput output

#### Example

**REQUIREMENT - Needs BAM file produced by hicup.py** Needs a .rmap file Needs a .baitmap file

When running the pipeline on a local machine without COMPSs:

```
1 python process_bam2chicago_tool.py \  
2     --config tests/json/config_bam2chicago.json \  
3     --in_metadata tests/json/input_bam2chicago.json \  
4     --out_metadata tests/json/output_bam2chicago.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_bam2chicago_tool.py \
7   --config tests/json/config_bam2chicago.json \
8   --in_metadata tests/json/input_bam2chicago.json \
9   --out_metadata tests/json/output_bam2chicago.json

```

### 3.5.2 Methods

**class** process\_bam2chicago\_tool.process\_bam2chicago (*configuration=None*)

This class creates .chinput files. input files for CHiCAGO

**run** (*input\_files, metadata, output\_files*)

This is the main function that run the tools to create .chinput files

#### Parameters

- **input\_files** (*dict*) –
  - BAM:** **str** path to BAM files
  - RMAP:** **str** path to RMAP file
  - BAITMAP:** **str** path to BAITMAP file
- **metadata** (*dict*) – input metadata

#### Returns

- *bool*
- **output\_metadata** (*dict*) – metadata for .chinput file

## 3.6 Data normalization and peak calling

This pipeline runs the normalization of the data and call the real chomatine interactions

### 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

**output\_dir:** directory with all output folders and files

## Example

### REQUIREMENT - Needs a reference genome

- Needs file with the capture sequences with FASTA format
  - settings file
  - **design dir:** .rmap .baitmap .npb .nbpb .poe

When running the pipeline on a local machine without COMPSs:

```
1 python process_run_chicago.py \  
2   --config tests/json/config_chicago.json \  
3   --in_metadata tests/json/input_chicago.json \  
4   --out_metadata tests/json/output_chicago.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_runChicago.py \  
7   --config tests/json/config_chicago.json \  
8   --in_metadata tests/json/input_chicago.json \  
9   --out_metadata tests/json/output_chicago.json
```

## 3.6.2 Methods

**class** process\_run\_chicago.process\_run\_chicago (configuration=None)

Function for processing capture Hi-C fastq files. Files are aligned, filtered and analysed for Capture Hi-C peaks

**run** (input\_files, metadata, output\_files)

This main function that run the chicago pipeline with runChicago.R wrapper

#### Parameters

- **input\_files** (*dict*) – location with the .chinput files. chinput\_file: str in case there is one input file chinput\_file: comma separated list in case there is more than one input file.
- **metadata** (*dict*) – Input metadata, str
- **output** (*dict*) – output file locations

#### Returns

- **output\_files** (*dict*) – Folder location with the output files
- **output\_metadata** (*dict*) – Output metadata for the associated files in output\_files

## 3.7 Run the entire CHI-C pipeline



---

Tools for processing fastq C-HiC files

---

## 4.1 Map and parser reads

### 4.1.1 hicup\_tool

**class** CHiC.tool.hicup\_tool.**hicup** (*configuration=None*)

Tool to run hicup, from fastq to bam files

**digest\_genome** (*genome\_name, re\_enzyme, genome\_loc, re\_enzyme2*)

This function takes a genome and digest it using a restriction enzyme specified

**Parameters**

- **genome\_name** (*str*) – name of the output genome
- **re\_enzyme** (*str*) – name of the enzyme used to cut the genome format example A<sup>^</sup>GATCT, BglII .
- **genome\_loc** (*str*) – location of the genome in FASTA format
- **re\_enzyme2** (*str*) – Restriction site 2 refers to the second, optional (other DNA shearing techniques such as sonication may be used) enzymatic digestion. This restriction site does NOT form a Hi-C ligation junction. This is the restriction enzyme that is used when the Hi-C sonication protocol is not followed. Typically the sonication protocol is followed.

**static get\_hicup\_params** (*params*)

Function to handle to extraction of commandline parameters and formatting them for use with hicup

**Parameters** **params** (*dict*) –

<b>--bowtie</b>	Specify the path to Bowtie
<b>--bowtie2</b>	Specify the path to Bowtie 2
<b>--config</b>	Specify the configuration file
<b>--digest</b>	Specify the digest file listing restriction fragment co-ordinates

<b>--example</b>	Produce an example configuration file
<b>--format</b>	Specify FASTQ format Options: Sanger, Solexa_Illumina_1.0, Illumina_1.3, Illumina_1.5
<b>--help</b>	Print help message and exit
<b>--index</b>	Path to the relevant reference genome Bowtie/Bowtie2 indices
<b>--keep</b>	Keep intermediate pipeline files
<b>--longest</b>	Maximum allowable insert size (bps)
<b>--nofill</b>	Hi-C protocol did NOT include a fill-in of sticky ends prior to ligation step and therefore FASTQ reads shall be truncated at the Hi-C restriction enzyme cut site (if present) sequence is encountered
<b>--outdir</b>	Directory to write output files
<b>--quiet</b>	Suppress progress reports (except warnings)
<b>--shortest</b>	Minimum allowable insert size (bps)
<b>--temp</b>	Write intermediate files (i.e. all except summaryfiles and files generated by HiCUP Deduplicator) to a specified directory
<b>--threads</b>	Specify the number of threads, allowing simultaneous processing of multiple files
<b>--version</b>	Print the program version and exit
<b>--zip</b>	Compress output

### Returns

**Return type** list

**hicup\_align\_filt** (*\*\*kwargs*)

This function aligns the HiC read into a reference genome and filter them

### Parameters

- **bowtie2\_loc** –
- **genome\_index** (*str*) – location of genome indexed with bowtie2
- **digest\_genome** (*str*) – location of genome digested
- **fastq1** (*str*) – location of fastq1 file
- **fastq2** (*str*) – location of fastq2

### Returns

**Return type** Bool

**run** (*input\_files, input\_metadata, output\_files*)

Function that runs and pass the parameters for all the functions

### Parameters

- **input\_files** (*dict*) –
- **metadata** (*dict*) –

- **output\_files** (*dict*) –

**untar\_index** (*\*\*kwargs*)

Extracts the Bowtie2 index files from the genome index tar file. :param genome\_file\_name: Location string of the genome fasta file :type genome\_file\_name: str :param genome\_idx: Location of the Bowtie2 index file :type genome\_idx: str :param bt2\_1\_file: Location of the <genome>.1.bt2 index file :type bt2\_1\_file: str :param bt2\_2\_file: Location of the <genome>.2.bt2 index file :type bt2\_2\_file: str :param bt2\_3\_file: Location of the <genome>.3.bt2 index file :type bt2\_3\_file: str :param bt2\_4\_file: Location of the <genome>.4.bt2 index file :type bt2\_4\_file: str :param bt2\_rev1\_file: Location of the <genome>.rev.1.bt2 index file :type bt2\_rev1\_file: str :param bt2\_rev2\_file: Location of the <genome>.rev.2.bt2 index file :type bt2\_rev2\_file: str

**Returns** Boolean indicating if the task was successful

**Return type** bool

## 4.2 Create CHiCAGO input files

### 4.2.1 makeRmap

### 4.2.2 makeBaitmap

### 4.2.3 makeDesignFiles

**class** CHiC.tool.makeDesignFiles.**makeDesignFilesTool** (*configuration=None*)

Tool for making the design files as part of the input for Chicago capture Hi-C

**static get\_design\_params** (*params*)

This function handle chicago parameters, selecting the given ones and passing to the command line.

**makeDesignFiles** (*\*\*kwargs*)

make the design files and store it in the specify design folder. It is a wrapper of makeDesignFiles.py

#### Parameters

- **designDir** (*str*,) – Path to the folder with the output files(recommended the same folder as .map and .baitmap files).
- **parameters** (*dict*,) – list of parameter already selected by get\_makeDesignFiles\_params().

#### Returns

- *bool*
- **outFilePrefix** (*str*) – writes the output files in the defined location

**run** (*input\_files, input\_metadata, output\_files*)

The main function to run makeDesignFiles.

#### Parameters

- **input\_files** (*dict*) – designDir : path to the designDir containin .rmap and .baitmap files
- **input\_metadata** (*dict*) –
- **output\_files** (*dict*) –

**outFilePrefix** [path to the output folder and prefix name of files] example: “/folder1/folder2/prefixname”. Recommended to use the path to designDir and the same prefix as .rmap and .baitmap

#### Returns

- **output\_files** (*dict*) – List of location for the output files.
- **output\_metadata** (*dict*) – List of matching metadata dict objects.

## 4.3 Convert bam files into chicago input

### 4.3.1 bam2chicago

**class** CHiC.tool.bam2chicago\_tool.bam2chicagoTool (*configuration=None*)

Tool for preprocess the input files

**bam2chicago** (*\*\*kwargs*)

Main function that preprocess the bam files into Chinput files. Part of the input files of CHiCAGO. It is a wrapper of bam2chicago.sh.

#### Parameters

- **bamFile** (*str*,) – path to paired-end file produced by a HiC aligner; Chicago has only been tested with data produced by HiCUP (<http://www.bioinformatics.babraham.ac.uk/projects/hicup/>). However, it should theoretically be possible to use other HiC aligners for this purpose.
- **rmapFile** (*str*,) – A tab-separated file of the format <chr> <start> <end> <numeric ID>, describing the restriction digest (or “virtual digest” if pooled fragments are used). These numeric IDs are referred to as “otherEndID” in Chicago. All fragments mapping outside of the digest coordinates will be disregarded by both these scripts and Chicago.
- **baitMapFile** (*str*,) – Tab-separated file of the format <chr> <start> <end> <numeric ID> <annotation>, listing the coordinates of the baited/captured restriction fragments (should be a subset of the fragments listed in rmapfile), their numeric IDs (should match those listed in rmapfile for the corresponding fragments) and their annotations (such as, for example, the names of baited promoters). The numeric IDs are referred to as “baitID” in Chicago.
- **chinput** (*str*) – name of the output file. Bbam2chicago creates a folder with the name of this sample, and inside the folder there is a file with chinput.chinput, that is the final output.

#### Returns

- *bool*
- **chinput** (*str*,) – name of the sample

**run** (*input\_files, input\_metadata, output\_files*)

Function that runs and pass the parameters to bam2chicago

#### Parameters

- **input\_files** (*dict*) –
- **hicup\_outdir\_tar** (*str*) –
- **rmapFile** (*str*) –

- **baitmapFile** (*str*) –
- **metadata** (*dict*) –

**Returns**

- **output\_files** (*list*)
- *List of locations for the output files.*
- **output\_metadata** (*list*)
- *List of matching metadata dict objects*

**sort\_chicago** (*\*\*kwargs*)

This function sort bamfile by name of the reads as bam2chicago requires

**Parameters**

- **bamfile** (*str*) –
- **sorted\_bam** (*str*) –

**Returns** sorted\_bam**Return type** str**static untar\_hicup\_out** (*hicup\_outdir\_tar, bam\_name*)

Untar hicup output folder

**Parameters**

- **hicup\_outdir\_tar** (*str*) – path to hicup output folder
- **path\_bam** (*str*) – path to bam file

**Returns****Return type** bool

## 4.4 Normalize data and call C-HiC peaks

### 4.4.1 run\_chicago

**class** CHiC.tool.run\_chicago.**ChicagoTool** (*configuration=None*)

tool for running the CHiCAGO algorithm

**chicago** (*\*\*kwargs*)

Run and annotate the Capture-HiC peaks. Chicago will create 4 folders under the output\_prefix data : output\_index.Rds -> chicago data saved on Rds format output\_index\_params.txt -> parameters used to run Chicago output\_index.export\_format -> chicago output in the chosen format diag\_plots : 3 plots to assess the quality of the output (see CHicago Capture-HiC documentation for details) enrichment\_data: files for the feature enrichment output (in case is used) examples: output\_index\_proxExamples.pdf: random chosen peaks showing interactions regions see <http://regulatorygenomicsgroup.org/chicago> for more information

**Parameters**

- **input\_files** (*str or comma separated list if there is more than one replicate*) –
- **output\_prefix** (*str*) –
- **output\_dir** (*str (whole path for the output)*) –

- **params** (*dict*) –

**Returns** writes the output files in the defined location

**Return type** bool

**static get\_chicago\_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 run chicago for peak calling. The input files are .chinput and are transformed from BAM files using bam2chicago.sh input files could be just one file or a comma separated files from more than one biological replicate. Technical replicates should be pooled to one .chinput

**Parameters**

- **input\_files** (*dict*) – list of .chinput files, or str with a single .chinput file
- **input\_metadata** (*dict*) –
- **output\_files** (*dict with the output path*) –

**Returns**

- **output\_files** (*Dict*) – List of locations for the output files,
- **output\_metadata** (*Dict*) – List of matching metadata dict objects

**static untar\_chinput** (*chinput\_tar*)

This function take as input the tar chinput

**Parameters** **chinput\_tar** (*str*) – path to the tar file, the tar files should have the same prefix name as the tar file

**Returns**

**Return type** list of untar files

#### **5.1 25-09-2018 handling\_chr\_header branch merge with master**

This rmap\_tool.py from this branch take the chromosome format from the used the reference genome and output a file with two columns, dictionary like with number of the chromosome and the name of the chromosome from the reference genome. example 1 chr1 2 chr2 3 chr3 ect. . .

This file is passed to the makeBaitmap.py script and generate the .batimap file with the corresponding chromosome name. This is necessary as the rtrees used in makeBaitmap.py needs an integer instead of “chr” or any other format.

#### **5.2 15-10-2018 mm\_mods\_for\_makebaitmaps branch merge with master**

This branch contains some modifications from Mark to solve issues with pyCOMPSs regarding makeBaitmap.py tool

#### **5.3 10-12-2018 creation of the branch VM\_CR1 (VM current release version)**

This Branch contains all the changes that to run the pipeline in the COMPSs VM from BSC cluster.





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.



## CHAPTER 7

---

### Indices and tables

---

- `genindex`
- `modindex`
- `search`



### p

`process_bam2chicago_tool`, [10](#)  
`process_design`, [8](#)  
`process_hicup`, [7](#)  
`process_run_chicago`, [11](#)

### t

`tool`, [13](#)





## B

bam2chicago() (CHiC.tool.bam2chicago\_tool.bam2chicagoTool method), 16

bam2chicagoTool (class in CHiC.tool.bam2chicago\_tool), 16

## C

chicago() (CHiC.tool.run\_chicago.ChicagoTool method), 17

ChicagoTool (class in CHiC.tool.run\_chicago), 17

## D

digest\_genome() (CHiC.tool.hicup\_tool.hicup method), 13

## G

get\_chicago\_params() (CHiC.tool.run\_chicago.ChicagoTool static method), 18

get\_design\_params() (CHiC.tool.makeDesignFiles.makeDesignFilesTool static method), 15

get\_hicup\_params() (CHiC.tool.hicup\_tool.hicup static method), 13

## H

hicup (class in CHiC.tool.hicup\_tool), 13

hicup\_alig\_filt() (CHiC.tool.hicup\_tool.hicup method), 14

## M

makeDesignFiles() (CHiC.tool.makeDesignFiles.makeDesignFilesTool method), 15

makeDesignFilesTool (class in CHiC.tool.makeDesignFiles), 15

## P

process\_bam2chicago (class in process\_bam2chicago\_tool), 11

process\_bam2chicago\_tool (module), 10

process\_design (class in process\_design), 9

process\_design (module), 8

process\_hicup (class in process\_hicup), 8

process\_hicup (module), 7

process\_run\_chicago (class in process\_run\_chicago), 12

process\_run\_chicago (module), 11

## R

run() (CHiC.tool.bam2chicago\_tool.bam2chicagoTool method), 16

run() (CHiC.tool.hicup\_tool.hicup method), 14

run() (CHiC.tool.makeDesignFiles.makeDesignFilesTool method), 15

run() (CHiC.tool.run\_chicago.ChicagoTool method), 18

run() (process\_bam2chicago\_tool.process\_bam2chicago method), 11

run() (process\_design.process\_design method), 9

run() (process\_hicup.process\_hicup method), 8

run() (process\_run\_chicago.process\_run\_chicago method), 12

## S

sort\_chicago() (CHiC.tool.bam2chicago\_tool.bam2chicagoTool method), 17

## T

tool (module), 13

## U

untar\_chinput() (CHiC.tool.run\_chicago.ChicagoTool static method), 18

untar\_hicup\_out() (CHiC.tool.bam2chicago\_tool.bam2chicagoTool static method), 17

untar\_index() (CHiC.tool.hicup\_tool.hicup method), 15