

---

# Micro-C Documentation

*Release 0.1*

**Dovetail**

**Aug 01, 2023**



**CONTENTS:**

<b>1</b>	<b>Overview</b>	<b>3</b>
1.1	Before you begin . . . . .	4
1.2	Pre-Alignment . . . . .	5
1.3	From fastq to final valid pairs bam file . . . . .	6
1.4	Library QC . . . . .	11
1.5	Library Complexity . . . . .	14
1.6	Generating Contact Matrix . . . . .	15
1.7	Micro-C Comparative Analyses . . . . .	20
1.8	Conformation Analysis . . . . .	22
1.9	Micro-C Data Sets . . . . .	24
1.10	Support . . . . .	25
<b>2</b>	<b>Indices and tables</b>	<b>27</b>

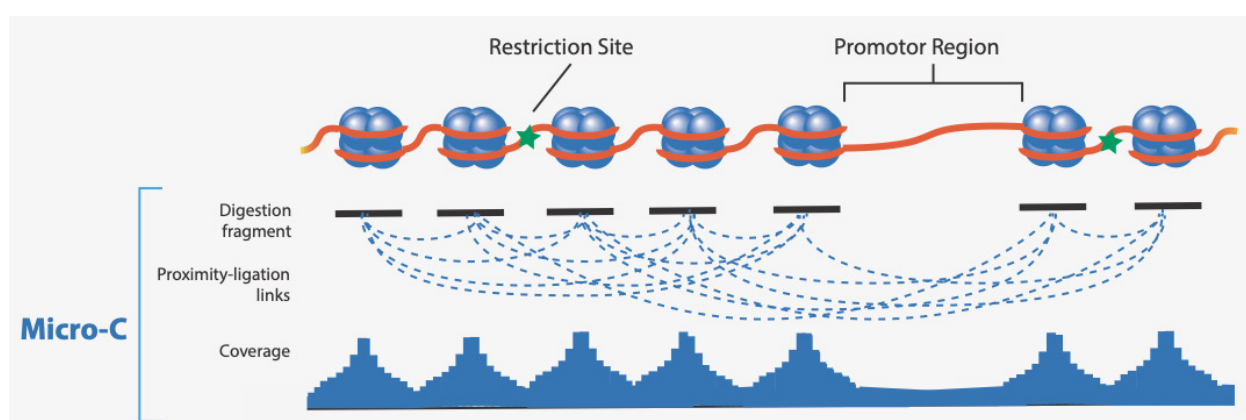






## OVERVIEW

- Dovetail™ Micro-C Kit uses the Micrococcal nuclease (MNase) enzyme instead of restriction enzymes for chromatin digestion, yielding 146 bp fragments distributed frequently across the genome.



- Key benefits of Micro-C:
  - Sequence-independent chromatin fragmentation enables even genome-wide detection of chromatin contacts (up to 20% of the genome lacks coverage using restriction enzyme based Hi-C approaches)
  - Ultra-high nucleosome-level resolution of chromatin contacts
  - Highest signal-to-noise data with both enrichment of long-range informative reads and nucleosome protected fragments
  - The ability to detect higher-order features, such as chromatin loops, in proximity ligation data is dependent on enriching long-range informative reads to capture chromatin interaction frequency. The increased chromosome conformation informative reads combined with ultra-high-resolution improves loop calling compared to RE-based methods.
- If you are using the Micro-C protocol as part of the Dovetail™ HiChIP MNase solution, please refer to our [HiChIP page](#) for further instructions.
- This guide will take you step by step on how to QC your Micro-C library, how to interpret the QC results and how to generate *contact maps*. If you don't yet have a sequenced Micro-C library and you want to get familiar with the data, you can download Micro-C sequences libraries from our publicly available [data sets](#).
- The QC process starts with aligning the reads to a reference genome then retaining high quality mapped reads. From there the mapped data will be used to generating a pairs file with pairtools, which categorizes pairs by read type and insert distance, this step both flags and removes PCR duplicates. Once pairs are categorized, counts of each class are summed and reported.
- If this is your first time following this tutorial, please check the [Before you begin page](#) first.

## 1.1 Before you begin

### 1.1.1 Have a copy of the Micro-C scripts on your machine:

Clone this repository:

```
git clone https://github.com/dovetail-genomics/Micro-C.git
```

### 1.1.2 Dependencies

Make sure that the following dependencies are installed:

- pysam
- tabulate
- bedtools
- deeptools
- matplotlib
- pandas
- bwa
- pairtools
- samtools
- preseq

If you are facing any issues with the installation of any of the dependencies, please contact the supporter of the relevant package.

python3 and pip3 are required, if you don't already have them installed, you will need sudo privileges.

- Update and install python3 and pip3:

```
sudo apt-get update
sudo apt-get install python3 python3-pip
```

- To set python3 and pip3 as primary alternative:

```
sudo update-alternatives --install /usr/bin/python python /usr/bin/python3 1
sudo update-alternatives --install /usr/bin/pip pip /usr/bin/pip3 1
```

If you are working on a new machine and don't have the dependencies, you can use the `installDep.sh` script in this repository for updating your instance and installing the dependencies and python3. This process will take approximately 10' and requires sudo privileges. The script was tested on Ubuntu 18.04 with the latest version as of 04/11/2020

If you choose to run the provided installation script you will first need to set the permission to the file:

```
chmod +x ./Micro-C/installDep.sh
```

And then run the installation script:

```
./Micro-C/installDep.sh
```



**Remember!**

Once the installation is completed, sign off and then sign back to your instance to refresh the database of applications.

### 1.1.3 Input files

For this tutorial you will need:

- **fastq files** R1 and R2, either fastq or fastq.gz are acceptable
- **reference in a fasta file format**, e.g. hg38

If you don't already have your own input files or want to run a test on a small data set, you can download sample fastq files from the [Micro-C Data Sets section](#). The 2M data set is suitable for a quick testing of the instruction in this tutorial.

```
wget https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_2M_R1.fastq
wget https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_2M_R2.fastq
```

## 1.2 Pre-Alignment

For downstream steps you will need a genome file, genome file is a tab delimited file with chromosome names and their respective sizes. If you don't already have a genome file follow these steps:

1. Generate an index file for your reference, a reference file with only the main chromosomes should be used (e.g. without alternative or unplaced chromosomes).

**Command:**

```
samtools faidx <ref.fasta>
```

**Example:**

```
samtools faidx hg38.fasta
```

Faidx will index the ref file and create <ref.fasta>.fai on the reference directory.

2. Use the index file to generate the genome file by printing the first two columns into a new file.

**Command:**

```
cut -f1,2 <ref.fasta.fai> > <ref.genome>
```

**Example:**

```
cut -f1,2 hg38.fasta.fai > hg38.genome
```

In line with the 4DN project guidelines and from our own experience optimal alignment results are obtained with Burrows-Wheeler Aligner (bwa). Prior to alignment, generate a bwa index file for the chosen reference.

```
bwa index <ref.fasta>
```

**Example:**

```
bwa index hg38.fasta
```

No need to specify an output path, the bwa index files are automatically generated at the reference directory. Please note that this step is time consuming, however you need to run it only once for a reference.

To avoid memory issues, some of the steps require writing temporary files into a temp folder, please generate a temp folder and remember its full path. Temp files may take up to x3 of the space that the fastq.gz files are taking, that is, if the total volume of the fastq files is 5Gb, make sure that the temp folder can store at least 15Gb.

**Command:**

```
mkdir <full_path/to/tmpdir>
```

**Example:**

```
mkdir /home/ubuntu/ebs/temp
```

In this example the folder *temp* will be generated on a mounted volume called *ebs* on a user account *ubuntu*.

## 1.3 From fastq to final valid pairs bam file

---

**fastq to final valid pairs bam file - for the impatient!**

If you just want to give it a shot and run all the alignment and filtering steps without going over all the details, we made a shorter version for you, with all the steps piped, outputting a final bam file with its index file and a dup stats file, otherwise move to the next section *fastq to final valid pairs bam file - step by step*

**Command:**

```
bwa mem -5SP -T0 -t<cores> <ref.fa> <MicroC.R1.fastq.gz> <MicroC.R2.fastq.gz> | \
pairtools parse --min-mapq 40 --walks-policy 5unique \
--max-inter-align-gap 30 --nproc-in <cores> --nproc-out <cores> --chroms-path <ref.
↳genome> | \
pairtools sort --tmpdir=<full_path/to/tmpdir> --nproc <cores>|pairtools dedup --nproc-in
↳<cores> \
--nproc-out <cores> --mark-dups --output-stats <stats.txt>|pairtools split --nproc-in
↳<cores> \
--nproc-out <cores> --output-pairs <mapped.pairs> --output-sam -|samtools view -bS -@
↳<cores> | \
samtools sort -@<cores> -o <mapped.PT.bam>;samtools index <mapped.PT.bam>
```

**Example:**

```
bwa mem -5SP -T0 -t16 hg38.fasta MicroC_2M_R1.fastq MicroC_2M_R2.fastq| pairtools parse -
↳--min-mapq 40 --walks-policy 5unique --max-inter-align-gap 30 --nproc-in 8 --nproc-out
↳8 --chroms-path hg38.genome | pairtools sort --tmpdir=/home/ubuntu/ebs/temp/ --nproc
↳16|pairtools dedup --nproc-in 8 --nproc-out 8 --mark-dups --output-stats stats.
↳txt|pairtools split --nproc-in 8 --nproc-out 8 --output-pairs mapped.pairs --output-
↳sam -|samtools view -bS -@16 | samtools sort -@16 -o mapped.PT.bam;samtools index
↳mapped.PT.bam
```



The full command above, with 2M read pairs on an Ubuntu 18.04 machine with 16 CPUs and 64GiB was completed in less than 5 minutes. On the same machine type.

### 1.3.1 fastq to final valid pairs bam file - step by step

#### Alignment

Now that you have a genome file, index file and a reference fasta file you are all set to align your Micro-C library to the reference. Please note the specific settings that are needed to map mates independently and for optimal results with our proximity library reads.

Parameter	Alignment function
mem	set the bwa to use the BWA-MEM algorithm, a fast and accurate alignment algorithm optimized for sequences in the range of 70bp to 1Mbp
-5	for split alignment, take the alignment with the smallest coordinate (5' end) as primary, the mapq assignment of the primary alignment is calculated independent of the 3' alignment
-S	skip mate rescue
-P	skip pairing; mate rescue performed unless -S also in use
-T0	The T flag set the minimum mapping quality of alignments to output, at this stage we want all the alignments to be recorded and thus T is set up to 0, (this will allow us to gather full stats of the library, at later stage we will filter the alignments by mapping quality)
-t	number of threads, default is 1. Set the numbers of threads to not more than the number of cores that you have on your machine (If you don't know the number of cores, used the command <code>lscpu</code> and multiply Thread(s) per core x Core(s) per socket x Socket(s))
*.fasta or *.fa	Path to a reference file, ending with .fa or .fasta, e.g, hg38.fasta
*.fastq or *.fastq.gz	Path to two fastq files; path to read 1 fastq file, followed by fastq file of read 2 (usually labeled as R1 and R2, respectively). Files can be in their compressed format (.fastq.gz) or uncompressed (.fastq). In case your library sequence is divided to multiple fastq files, you can use a process substitution < with the cat command (see example below)
-o	sam file name to use for output results [stdout]. You can choose to skip the -o flag if you are piping the output to the next command using ' '

Bwa mem will output a sam file that you can either pipe or save to a path using -o option, as in the example below (please note that version 0.7.17 or higher should be used, older versions do not support the -5 flag)

#### Command:

```
bwa mem -5SP -T0 -t<threads> <ref.fasta> <MicroC_R1.fastq> <MicroC_R2.fastq> -o <aligned.
↪sam>
```

#### Example (one pair of fastq files):

```
bwa mem -5SP -T0 -t16 hg38.fasta MicroC_2M_R1.fastq MicroC_2M_R2.fastq -o aligned.sam
```

#### Example (multiple pairs of fastq files):

```
bwa mem -5SP -T0 -t16 hg38.fasta <(zcat file1.R1.fastq.gz file2.R1.fastq.gz file3.R1.
↪fastq.gz)> <(zcat file1.R2.fastq.gz file2.R2.fastq.gz file3.R2.fastq.gz)> -o aligned.sam
```

## Recording valid ligation events

We use the `parse` module of the `pairtools` pipeline to find ligation junctions in Micro-C (and other proximity ligation) libraries. When a ligation event is identified in the alignment file the `pairtools` pipeline will record the outer-most (5') aligned base pair and the strand of each one of the paired reads into `.pairsam` file (pairsam format captures SAM entries together with the Hi-C pair information). In addition, it will also assign a pair type for each event. e.g. if both reads aligned uniquely to only one region in the genome, the type UU (Unique-Unique) will be assigned to the pair. The following steps are necessary to identify the high quality valid pairs over low quality events (e.g. due to low mapping quality):

`pairtools parse` options:

Parameter	Value	Function
<code>min-mapq</code>	40	Mapq threshold for defining an alignment as a multi-mapping alignment. Alignment with mapq <40 will be marked as type M (multi)
<code>walks-policy</code>	5unique	Walks is the term used to describe multiple ligations events, resulting three alignments (instead of two) for a read pair. However, there are cases in which three alignment in read pairs are the result of one ligation event, pairtool parse can rescue this event. walks-policy is the policy for reporting un-rescuable walk. 5unique is used to report the 5'-most unique alignment on each side, if present (one or both sides may map to different locations on the genome, producing more than two alignments per DNA molecule)
<code>max-inter-align-gap</code>	30	In cases where there is a gap between alignments, if the gap is 30 or smaller, ignore the gap, if the gap is >30bp, mark as "null" alignment
<code>nproc-in</code>	integer, e.g. 16	pairtools has an automatic-guess function to identify the format of the input file, whether it is compressed or not. When needed, the input is decompressed by <code>bgzip/lz4c</code> . The option <code>nproc-in</code> set the number of processes used by the auto-guessed input decompressing command, if not specified, default is 3
<code>nproc-out</code>	integer, e.g. 16	pairtools automatic-guess the desired format of the output file (compressed or not compressed, based on file name extension). When needed, the output is compressed by <code>bgzip/lz4c</code> . The option <code>nproc-out</code> set the number of processes used by the auto-guessed output compressing command, if not specified, default is 8
<code>chroms-path</code>		path to <code>.genome</code> file, e.g. <code>hg38.genome</code>
<code>*.sam</code>		path to sam file used as an input. If you are piping the input (stdin) skip this option
<code>*pairsam</code>		name of pairsam file for writing output results. You can choose to skip and pipe the output directly to the next command ( <code>pairtools sort</code> )

`pairtools parse` command example for finding ligation events:

### Command:

```
pairtools parse --min-mapq 40 --walks-policy 5unique --max-inter-align-gap 30 --nproc-in
↪ <cores> \
--nproc-out <cores> --chroms-path <ref.genome> <aligned.sam> > <parsed.pairsam>
```

### Example:

```
pairtools parse --min-mapq 40 --walks-policy 5unique --max-inter-align-gap 30 --nproc-in ↪
↪ 8 --nproc-out 8 --chroms-path hg38.genome aligned.sam > parsed.pairsam
```

At the parsing step, pairs will be flipped such that regardless of read1 and read2, pairs are always recorded with first side of the pair having the lower genomic coordinates.

## Sorting the pairsam file

The parsed pairs are then sorted using *pairtools sort*

`pairtools sort` options:

Parameter	Function
<code>--tmpdir</code>	Provide a full path to a temp directory. A good rule of thumb is to have a space available for this directory at a volume of x3 of the overall volume of the fastq.gz files. Using a temp directory will help avoid memory issues
<code>--nproc</code>	Number of processes to split the sorting work

### Command:

```
pairtools sort --nproc <cores> --tmpdir=<path/to/tmpdir> <parsed.pairsam> > <sorted.
↪pairsam>
```

### Example:

```
pairtools sort --nproc 16 --tmpdir=/home/ubuntu/ebs/temp/ parsed.pairsam > sorted.
↪pairsam
```

### Important!

Please note that an absolute path for the temp directory is required for `pairtools sort`, e.g. path of the structure `~/ebs/temp/` or `./temp/` will not work, instead, something of this sort is needed `/home/user/ebs/temp/`

## Removig PCR duplicates

`pairtools dedup` detects molecules that could be formed via PCR duplication and tags them as “DD” pair type. These pairs should be excluded from downstream analysis. Use the `pairtools dedup` command with the `--output-stats` option to save the dup stats into a text file.

`pairtools dedup` options:

Parameter	Function
<code>--mark-dups</code>	If specified, duplicate pairs are marked as DD in “pair_type” and as a duplicate in the sam entries
<code>--output-stats</code>	Output file for duplicate statistics. Please note that if a file with the same name already exists, it will be opened in the append mode

### Command:

```
pairtools dedup --nproc-in <cores> --nproc-out <cores> --mark-dups --output-stats <stats.
↪txt> \
--output <dedup.pairsam> <sorted.pairsam>
```

### Example:

```
pairtools dedup --nproc-in 8 --nproc-out 8 --mark-dups --output-stats stats.txt --output-  
↪ dedup.pairsam sorted.pairsam
```

## Generate .pairs and bam files

The `pairtools split` command is used to split the final `.pairsam` into two files: `.sam` (or `.bam`) and `.pairs` (`.pairsam` has two extra columns containing the alignments from which the Micro-C pair was extracted, these two columns are not included in `.pairs` files)

`pairtools split` options:

Parameter	Function
<code>--output-pairs</code>	Output pairs file. If the path ends with <code>.gz</code> or <code>.lz4</code> the output is <code>pbgzip/lz4c</code> -compressed. If you wish to pipe the command and output the pairs files to stdout use <code>-</code> instead of file name
<code>--output-sam</code>	Output sam file. If the file name extension is <code>.bam</code> , the output will be written in bam format. If you wish to pipe the command, use <code>-</code> instead of a file name. please note that in this case the sam format will be used (and can be later converted to bam file e.g. with the command <code>samtools view -bS -@16 -o temp.bam</code> )

### Command:

```
pairtools split --nproc-in <cores> --nproc-out <cores> --output-pairs <mapped.pairs> \  
--output-sam <unsorted.bam> <dedup.pairsam>
```

### Example:

```
pairtools split --nproc-in 8 --nproc-out 8 --output-pairs mapped.pairs --output-sam-  
↪ unsorted.bam dedup.pairsam
```

The `.pairs` file can be used for generating *contact matrix*

## Generating the final bam file

For downstream steps, the bam file should be sorted, using the command `samtools sort`

`samtools sort` options:

Parameter	Function
<code>-@</code>	number of threads to use
<code>-o</code>	file name. Write final output to FILE rather than standard output
<code>-T</code>	path to temp file. Using a temp file will help avoiding memory issues

### Command:

```
samtools sort -@<threads> -T <path/to/tmpdir/tempfile.bam>-o <mapped.PT.bam> <unsorted.  
↪ bam>
```

### Example:

```
samtools sort -@16 -T /home/ubuntu/ebs/temp/temp.bam -o mapped.PT.bam unsorted.bam
```

For future steps an index (.bai) of the bam file is also needed. Index the bam file:

**Command:**

```
samtools index <mapped.PT.bam>
```

**Example:**

```
samtools index mapped.PT.bam
```

The mapped.PT.bam is the final bam file that will be used downstream steps.

The above steps resulted in multiple intermediate files, to simplify the process and avoid intermediate files, you can pipe the steps as in the example above (*fastq to final valid pairs bam file - for the impatient*)

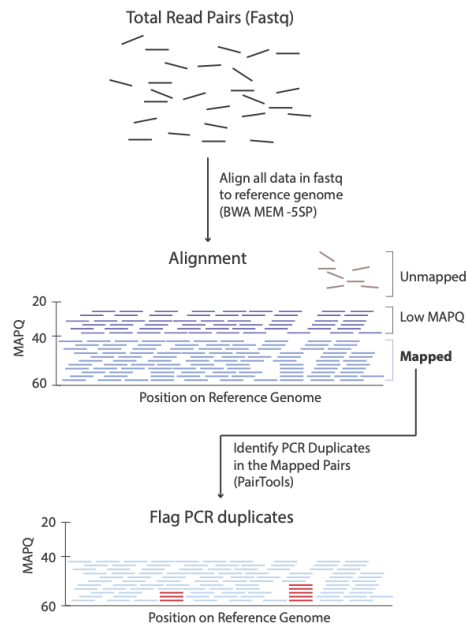
## 1.4 Library QC

At step *Removing PCR duplicates* you used the flag `-output-stats`, generating a stats file in addition to the pairsam output (e.g. `-output-stats stats.txt`). The stats file is an extensive output of pairs statistics as calculated by pairtools, including total reads, total mapped, total dups, total pairs for each pair of chromosomes etc'. Although you can use directly the pairtools stats file as is to get informed on the quality of the Micro-C library, we find it easier to focus on a few key metrics. We include in this repository the script `get_qc.py` that summarize the paired-tools stats file and present them in percentage values in addition to absolute values.

The images below explains how the values on the QC report are calculated:

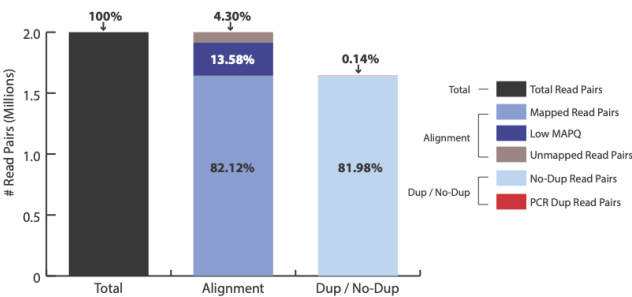
I. Aligning and filtering to remove low mapping quality and PCR duplicate read pairs

Process



Results

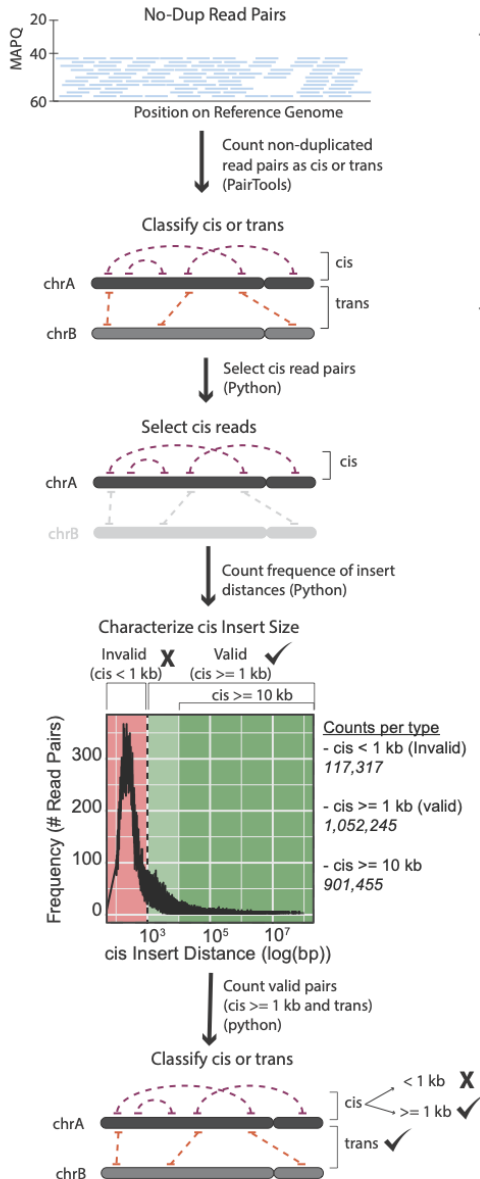
Category	Count	Percent
Total Read Pairs	2,000,000	100.00%
Unmapped Read Pairs	86,000	4.30%
Mapped Read Pairs	1,642,400	82.12%
PCR Dup Read Pairs	2,800	0.14%
No-Dup Read Pairs	1,640,800	81.98%
No-Dup Cis Read Pairs	1,169,563	71.28%
No-Dup Trans Read Pairs	471,317	28.72%
No-Dup Valid Read Pairs (cis >= 1 kb + trans)	1,523,483	92.85%
No-Dup Cis Read Pairs < 1kb	117,317	7.15%
No-Dup Cis Read Pairs >= 1kb	1,052,245	64.13%
No-Dup Cis Read Pairs >=10kb	901,455	54.94%





## II. Classifying read pairs (cis or trans), characterizing insert size, and identifying valid pairs

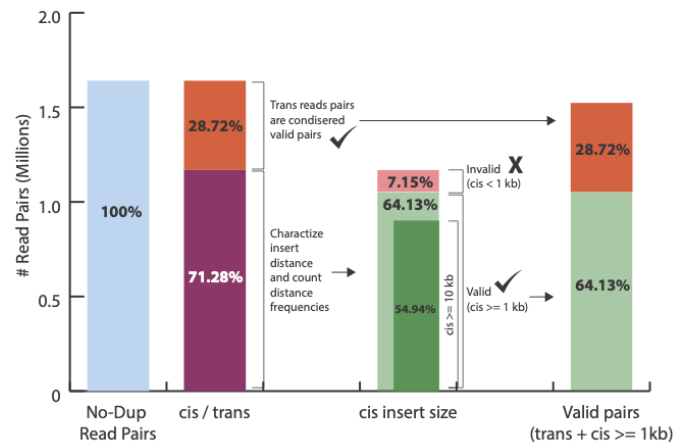
### Process



### Results

Category	Count	Percent
Total Read Pairs	2,000,000	100.00%
Unmapped Read Pairs	86,000	4.30%
Mapped Read Pairs	1,642,400	81.98%
PCR Dup Read Pairs	2,800	0.14%
No-Dup Read Pairs	1,640,800	82.04%
No-Dup Cis Read Pairs	1,169,563	71.28%
No-Dup Trans Read Pairs	471,317	28.72%
No-Dup Valid Read Pairs (cis >= 1 kb + trans)	1,523,483	92.85%
No-Dup Cis Read Pairs < 1 kb	117,317	7.15%
No-Dup Cis Read Pairs >= 1 kb	1,052,245	64.13%
No-Dup Cis Read Pairs >= 10kb	901,455	54.94%

Proportion of No-Dup Read Pairs



- No-Dup Read Pairs
- No-Dup Cis Read Pairs
- No-Dup Trans Read Pairs
- No-Dup Cis Read Pairs >= 1kb
- No-Dup Cis Read Pairs >= 10kb
- No-Dup Cis Read Pairs < 1kb
- ✓ Valid Read Pair
- ✗ Invalid Read Pair

Command:

```
python3 ./Micro-C/get_qc.py -p <stats.txt>
```

**Example:**

```
python3 ./Micro-C/get_qc.py -p stats.txt
```

After the script completes, it will print:

Total Read Pairs	2,000,000	100%
Unmapped Read Pairs	92,059	4.6%
Mapped Read Pairs	1,637,655	81.88%
PCR Dup Read Pairs	5,426	0.27%
No-Dup Read Pairs	1,632,229	81.61%
No-Dup Cis Read Pairs	1,288,943	78.97%
No-Dup Trans Read Pairs	343,286	21.03%
No-Dup Valid Read Pairs (cis >= 1kb + trans)	1,482,597	90.83%
No-Dup Cis Read Pairs < 1kb	149,632	9.17%
No-Dup Cis Read Pairs >= 1kb	1,139,311	69.8%
No-Dup Cis Read Pairs >= 10kb	870,490	53.33%

## 1.5 Library Complexity

If you performed a shallow sequencing experiment (e.g. 2M reads) and running a QC analysis to decide which library to use for deep sequencing (DS), it is recommended to evaluate the complexity of the library before moving to DS.

The *lc\_extrap* utility of the *preseq* package aims to predict the complexity of sequencing libraries.

*preseq* options:

Parameter	Value	Function
bam		Specifies that the input file type is bam. Please note that for a bam file to be a recognized input file htlib should be installed as well and preseq should be built with htlib support (for more details see preseq documentation or our installDep.sh script as example)
pe		Specifies that paired end data is used
extrap	2.10E+09	Maximum extrapolation
step	1.00E+08	Extrapolation step size
seg_len	1000000000	maximum segment length when merging paired end bam
output		output file

Please note that the input bam file should be a version prior to dups removal.

*preseq lc\_extrap* command example for extrapolating library complexity:

**Command:**

```
preseq lc_extrap -bam -pe -extrap 2.1e9 -step 1e8 -seg_len 1000000000 -output <output_
↪file> <input bam file>
```

**Example:**

```
preseq lc_extrap -bam -pe -extrap 2.1e9 -step 1e8 -seg_len 1000000000 -output out.preseq_
↪mapped.PT.bam
```

In this example the output file *out.preseq* will detail the extrapolated complexity curve of your library, with the number of reads in the first column and the expected distinct read value in the second column. For a typical experiment (human sample) check the expected complexity at 300M reads (to show the content of the file, type **cat out.preseq**). Expected unique pairs at 300M sequencing is at least ~ 120 million.

TOTAL_READS	EXPECTED_DISTINCT	LOWER_0.95CI	UPPER_0.95CI
0 0	0 0		
100000000.0	87353657.3	87094889.8	87570833.4
200000000.0	156073911.7	155222525.2	156597384.7
300M 300000000.0	211566844.9	209951454.2	212512486.5
400000000.0	257305075.8	254854778.3	258725614.9
500000000.0	295664757.3	292325778.7	297568930.4
600000000.0	328296017.0	324091795.9	330675142.4
700000000.0	356347398.8	351363600.4	359222765.3
800000000.0	380763483.3	375031977.1	384092719.4
900000000.0	402215465.5	395766880.5	405952677.6
1000000000.0	421214292.2	414081873.4	425317912.5
1100000000.0	438141525.9	430377245.4	442592433.2
1200000000.0	453314818.1	444969611.9	458097479.1
1300000000.0	466999393.9	458112691.8	472091684.1
1400000000.0	479404104.3	470012187.6	484785619.1
1500000000.0	490700479.3	480836631.0	496352461.5
1600000000.0	501030708.0	490725416.5	506935942.7
1700000000.0	510513655.0	499794837.4	516656351.6
1800000000.0	519249454.4	508142681.2	525615126.4
1900000000.0	527323059.6	515851769.5	533898408.5
2000000000.0	534807014.9	522992716.2	541579821.0

## 1.6 Generating Contact Matrix

There are two common formats for contact maps, the [Cooler format](#) and [Hic format](#). Both are compressed and sparsed formats to avoid large storage volumes; For a given  $n$  number of bins in the genome, the size of the matrix would be  $n^2$ , in addition, typically more than one resolution (bin size) is being used.

In this section we will guide you on how to generate both matrices types, *HiC* and *cool* based on the *.pairs file* that you generated in the [previous section](#) and how to visualize them.

## 1.6.1 Generating HiC contact maps using Juicer tools

### Additional Dependencies

- **Juicer Tools** - Download the JAR file for juicertools and place it in the same directory as this repository and name it as `juicertools.jar`. You can find the link to the most recent version of Juicer tools [here](#) e.g.:

```
wget https://s3.amazonaws.com/hicfiles.tc4ga.com/public/juicer/juicer_tools_1.22.01.jar
mv juicer_tools_1.22.01.jar ./Micro-C/juicertools.jar
```

- **Java** - If not already installed, you can install Java as follows:

```
sudo apt install default-jre
```

### From .pairs to .hic contact matrix

- **Juicer Tools** is used to convert `.pairs` file into a **HiC** contact matrix.
- HiC is highly compressed binary representation of the contact matrix
- Provides rapid random access to any genomic region matrix
- Stores contact matrix at 9 different resolutions (2.5M, 1M, 500K, 250K, 100K, 50K, 25K, 10K, and 5K)
- Can be programmatically manipulated using straw python API

The `.pairs` file that you generated in the *From fastq to final valid pairs bam file* section can be used directly with Juicer tools to generate the *HiC* contact matrix:

Parameter	Function
-Xmx	The flag Xmx specifies the maximum memory allocation pool for a Java virtual machine, from our experience 48000m works well when processing human data sets, If you are not sure how much memory your system has, run the command <code>free -h</code> and check the value under <i>total</i> .
Djava.awt.headless=true	Java is ran in a headless mode when the application does not interact with a user (if not specified, the default is <code>Djava.awt.headless=false</code> )
pre	The pre command allows users to create <code>.hic</code> files from their own data
--threads	Specifies the numbers of threads to be used (integer number)
*.pairs or *.pairs.gz	input file for generating the contact matrix
*.genome	genome file, listing the chromosomes and their sizes
*.hic	hic output file, containing the contact matrix

---

#### Tip no.1

Please note that if you have an older version of **Juicer tools**, generating contact map directly from `.pairs` file may not be supported. We recommend updating to a newer version. As we tested, the `pre` utility of the version 1.22.01 support the `.pairs` to HiC function.

---

#### Command:

```
java -Xmx<memory> -Djava.awt.headless=true -jar <path_to_juicer_tools.jar> pre --
↳ threads <no_of_threads> <mapped.pairs> <contact-map.hic> <ref.genome>
```

#### Example:

```
java -Xmx48000m -Djava.awt.headless=true -jar ./Micro-C/juicertools.jar pre --threads_
↳16 mapped.pairs contact_map.hic hg38.genome
```

---

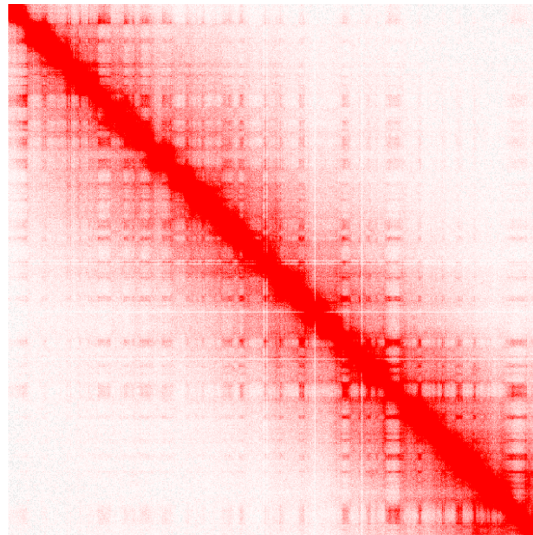
**Tip no.2**

Juicer tools offers additional functions that were not discussed here, including matrix normalization and generating matrix for only specified regions in the genome. To learn more about advanced options, please refer to the [Juicer Tools documentation](#).

---

**Visualizing .hic contact matrix**

The visualization tool Juicebox can be used to visualize the contact matrix. You can either [download](#) a local version of the tool to your computer as a Java application or use a [web](#) version of Juicebox. Load your .hic file to visualize the contact map and zoom in to areas of interest.



## 1.6.2 Generating cooler contact maps

**Additional Dependencies****Installing Cooler and its dependencies**

- `libhdf5 - sudo apt-get install libhdf5-dev`
- `h5py - pip3 install h5py`
- `cooler - pip3 install cooler`

For any issues with cooler installation or its dependencies, please refer to the [cooler installation documentation](#)

## Installing Pairix

**Pairix** is a tool for indexing and querying on a block-compressed text file containing pairs of genomic coordinates. You can install it directly from its github repository as follows:

```
git clone https://github.com/4dn-dcic/pairix
cd pairix
make
```

Add the bin path, and utils path to PATH and exit the folder:

```
PATH=~/.pairix/bin/:~/.pairix/util:~/.pairix/bin/pairix:$PATH
cd ..
```

---

### Important!

make sure to modify the following example with the path to your *pairix* installation folder. If you are not sure what is the path you can check it with the command *pwd* when located in the *pairix* folder.

---

For any issues with *pairix*, please refer to the [pairix documentation](#)

## From .pairs to cooler contact matrix

- **Cooler tools** is used to convert **indexed .pairs** file into **cool** and **mcool** contact matrices
- Cooler generates a sparse, compressed, and binary persistent representation of proximity ligation contact matrix
- Store matrix as **HDF5** file object
- Provides python API to manipulate contact matrix
- Each cooler matrix is computed at a specific resolution
- Multi-cool (mcool) files store a set of cooler files into a single HDF5 file object
- Multi-cool files are helpful for visualization

## Indexing the .pairs file

We will use the `cload pairix` utility of Cooler to generate contact maps. This utility requires the `.pairs` file to be indexed. **Pairix** is used for indexing compressed `.pairs` files. The files should be compressed with **bgzip** (which should already be installed on your machine). If your `.pairs` file is not yet **bgzip** compressed, first compress it as follows:

### Command:

```
bgzip <mapped.pairs>
```

### Example:

```
bgzip mapped.pairs
```

Following this command `mapped.pairs` will be replaced with its compressed form `mapped.pairs.gz`

---

### Note!

Compressing the `.pairs` file with `gzip` instead of `bgzip` will also result in a compressed file with the `.gz` suffix, but due to format differences it will not be accepted as an input for `pairix`.

Next, index the file `.pairs.gz` file:

**Command:**

```
pairix <mapped.pairs.gz>
```

**Example:**

```
pairix mapped.pairs.gz
```

## Generating single resolution contact map files

As mentioned above, we will use the `cloud pairix` utility of Cooler to generate contact maps:

`cooler cloud pairix` usage:

Parameter	Function
<genome_files>:<bin size>	Specifies the reference <i>.genome file</i> , followed with ``:`` and the desired bin size in bp
-p	Number of processes to split the work between (integer), default: 8
*.pairs.gz	Path to bgzip compressed and indexed <code>.pairs</code> file
*.cool	Name of output file

**Command:**

```
cooler cloud pairix -p <cores> <ref.genome>:<bin_size_in_bp> <mapped.pairs.gz> <matrix.  
→cool>
```

**Example:**

```
cooler cloud pairix -p 16 hg38.genome:1000 mapped.pairs.gz matrix_1kb.cool
```

## Generating multi-resolutions files and visualizing the contact matrix

When you wish to visualize the contact matrix, it is highly recommended to generate a multi-resolution `.mcool` file to allow zooming in and out to inspect regions of interest. The `cooler zoomify` utility allows you to generate a multi-resolution cooler file by coarsening. The input to `cooler zoomify` is a single resolution `.cool` file, to allow zooming in into regions of interest we suggest to generate a `.cool` file with a small bin size, e.g. 1kb. Multi-resolution files uses the suffix `.mcool`.

`cooler zoomify` usage:

Parameter	Function
-balance	Apply balancing to each zoom level. Off by default
-p	Number of processes to use for batch processing chunks of pixels, default: 1
*.cool	Name of contact matrix input file

*Command:\**

```
cooler zoomify --balance -p <cores> <matrix.cool>
```

**Example:**

```
cooler zoomify --balance -p 16 matrix_1kb.cool
```

The example above will result in a new file named *matrix\_1kb.mcool* (no need to specify output name)

---

**Tip**

Cooler offers additional functions that were not discussed here, including generating a cooler from a pre-binned matrix, matrix normalization and more. To learn more about advanced options, please refer to the cooler [documentation](#)

---

HiGlass is an interactive tool for visualizing .mcool files. To learn more about how to set up and use HiGlass follow the HiGlass [tutorial](#)

## 1.7 Micro-C Comparative Analyses

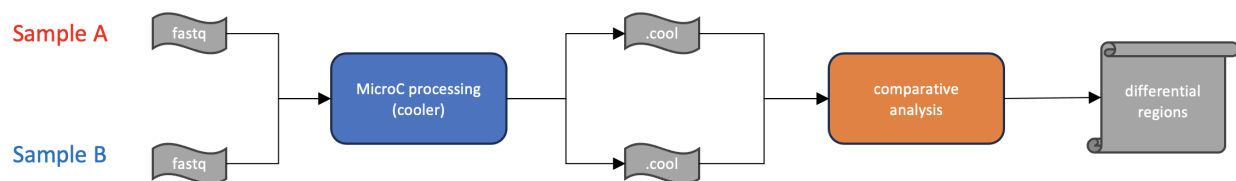
### 1.7.1 Introduction

Biological questions are seldom answered by analysing single samples in isolation. It is often the case that an experiment aims to make comparisons between two (or more) biological conditions, such as:

- 1) Untreated wild type vs treatment
- 2) Wild type vs knockout
- 3) Normal sample vs tumor

In all cases the goal is to produce a list of differentially interacting regions in one condition relative to the other. The main output for comparative analyses is analogous to what is expected for differential gene expression, where the primary result is a table of regions, the fold change between conditions, and a statistical measure of significance. For Micro-C, we aim to identify regions of differential interaction directly from the matrix files. See [previous steps](#) to generate the required matrices for differential analysis.

Figure 1:





### 1.7.2 Differential Analysis

**Question:** How do I perform differential analyses for Micro-C experiments?

**Process:** Mcool files are first converted to text files of a preferred resolution, and then used as input to the HiCcompare algorithm.

**Results:** Final results consist of a table of differentially interacting regions, fold change, and measure of statistical significance.

**Files and tools needed:**

- .cool, .mcool, .hic, or Hic-Pro files for each replicate and sample condition
- [HiCcompare](#) for single-replicate analysis or [multiHiCcompare](#) for multiple replicate experiments.

As the design of differential analysis experiments are unique to each biological question, there are multiple possibilities for how the analysis can be set up. A common scenario is to compare two conditions where each condition has two replicates, and is described in the [multiHiCcompare vignette](#). The HiCcompare package also contains functions for conversion of various input files

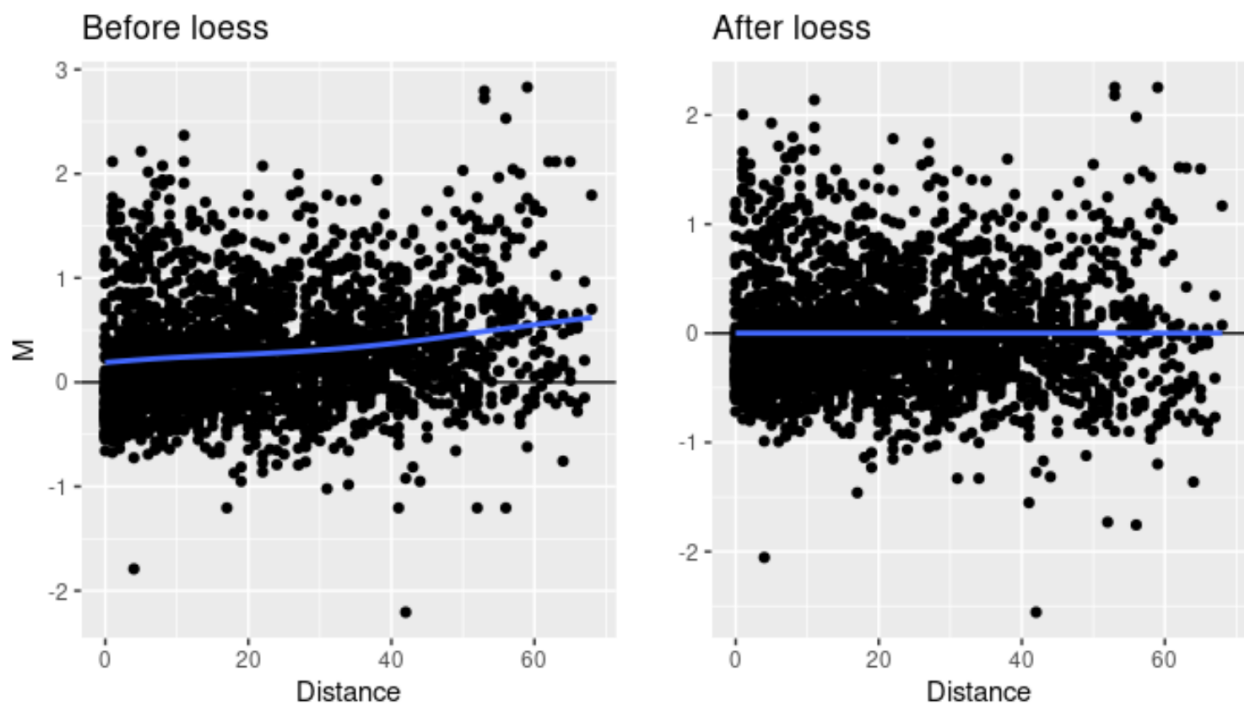
**Interpreting results:**

Micro-C differential analysis produces a number of intermediate files in addition to the final results table. There are two main outputs to consider:

- 1) MD normalization plots
- 2) Differential regions table

MD is a concept introduced by the HiCcompare developers and is analogous to the Tukey's mean/difference plot. M corresponds to the log2 fold change between the two conditions, and D is the distance between the two interacting regions. Loess normalization aims to eliminate the bias introduced by the influence of interaction distance on fold change between two conditions. It is often useful to visualize the effect of normalization between conditions to ensure the data is appropriate for downstream difference detection. An example effect of normalization is given below:

Figure 2:



For difference detection, the resulting output file is highly similar to what is expected for gene expression studies, where regions are listed and prioritized by a combination of fold change and a measure of statistical significance. Below is an example output from HicCompare:

chr1	start1	end1	chr2	start2	end2	IF1	IF2	D	M	adj.IF	adj.IF2	adj.M	mc	A	Z	p.val	p.adj
chr1	10000	11000	chr1	10000	11000	15	1	0	-	14.207	1.056	-	-	7.631	-	0.000	0.736
									3.907			3.750	0.157		3.603		
chr1	16000	17000	chr1	16000	17000	6	2	0	-	5.683	2.112	-	-	3.897	-	0.197	0.863
									1.585			1.428	0.157		1.291		
chr1	17000	18000	chr1	17000	18000	6	3	0	-	5.683	3.167	-	-	4.425	-	0.479	0.904
									1.000			0.843	0.157		0.708		
chr1	22000	23000	chr1	22000	23000	3	1	0	-	2.841	1.056	-	-	1.949	NA	1.000	1.000
									1.585			1.428	0.157				
chr1	24000	25000	chr1	24000	25000	1	1	0	0.000	0.947	1.056	0.157	-	1.001	NA	1.000	1.000
													0.157				
chr1	27000	28000	chr1	27000	28000	2	2	0	0.000	1.894	2.112	0.157	-	2.003	0.288	0.773	0.904
													0.157				
chr1	28000	29000	chr1	28000	29000	1	1	0	0.000	0.947	1.056	0.157	-	1.001	NA	1.000	1.000
													0.157				
chr1	31000	32000	chr1	31000	32000	4	1	0	-	3.788	1.056	-	-	2.422	-	0.088	0.863
									2.000			1.843	0.157		1.704		
chr1	36000	37000	chr1	36000	37000	2	1	0	-	1.894	1.056	-	-	1.475	NA	1.000	1.000
									1.000			0.843	0.157				

The most relevant fields from the output will be:

- adj.M – the log fold change in coverage between the two conditions
- p.adj – a p-value, after correction for multiple hypothesis testing, on the statistical significance of the observed fold change

#### Considerations:

- Replication – It is generally advisable to have technical replicates for differential analyses, as this will produce more statistically robust results.

## 1.8 Conformation Analysis

There are many open-source tools available that enable researchers to perform conformation analyses on contact matrices. We aim to highlight the tools and commands we use here at Cantata Bio call features. This is not a comprehensive list, nor do we own or manage any of the tools listed below. Please refer to their source document pages for more information or for help in trouble shooting the use of these tools. For each analysis performed below we provide, a link to the tool repo, the input file, the example command, and an example output for you to check your results against.

### 1.8.1 Example input files (.mcool and .hic files)

- Example mcool file
- Example hic file

### 1.8.2 A/B Compartments

#### Recommended Software

- Fanc-C <https://fan-c.readthedocs.io/en/latest/>

#### Example Command

```
fanc compartments -f -v MicroC_800M_eigen_64kb.bed -d MicroC_800M_64kb.bed -g hg38.fa
↪MicroC_800M.mcool@64000 MicroC_800M_64kb.ab
```

#### Example Output(s)

- Example AB compartments file
- Example AB compartments bed file
- Example AB Eigenvector file

### 1.8.3 Topologically Associated Domains

#### Recommended Software

- Juicer Arrowhead <https://github.com/aidenlab/juicer>

#### Example Command

```
java -jar -Xmx48000m -Djava.awt.headless=true -jar juicer_tools.jar arrowhead --threads
↪16 -k KR -m 2000 -r 10000 MicroC_800M.hic TAD_calls
java -jar -Xmx48000m -Djava.awt.headless=true -jar juicer_tools.jar arrowhead --threads
↪16 -k KR -m 2000 -r 5000 MicroC_800M.hic TAD_calls
```

#### Example Output(s)

- Example 10kb TADs file
- Example 5kb TADs file

### 1.8.4 Chromatin Loops

#### Recommended Software

- Mustache <https://github.com/ay-lab/mustache>

#### Example Command

```
mustache -p 48 -f MicroC_800M.mcool -r 16000 -o MicroC_800M_16000kb_loops.tsv
mustache -p 48 -f MicroC_800M.mcool -r 4000 -o MicroC_800M_4000kb_loops.tsv
```

#### Example Output(s)

- Example 16kb Loops file

- Example 4kb Loops file

## 1.9 Micro-C Data Sets

To download one of the data sets, simply use the wget command:

```
wget https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_100M_R1.fastq
wget https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_100M_R2.fastq
```

For testing purposes, we recommend using the 2M reads data sets, for any other purpose we recommend using the 800M reads data set.

Library	Link
GM12878 Micro-C 2M	<ul style="list-style-type: none"><li>• <a href="https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_2M_R1.fastq">https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_2M_R1.fastq</a></li><li>• <a href="https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_2M_R2.fastq">https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_2M_R2.fastq</a></li></ul>
GM12878 Micro-C 100M	<ul style="list-style-type: none"><li>• <a href="https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_100M_R1.fastq">https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_100M_R1.fastq</a></li><li>• <a href="https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_100M_R2.fastq">https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_100M_R2.fastq</a></li></ul>
GM12878 Micro-C 200M	<ul style="list-style-type: none"><li>• <a href="https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_200M_R1.fastq">https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_200M_R1.fastq</a></li><li>• <a href="https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_200M_R2.fastq">https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_200M_R2.fastq</a></li></ul>
GM12878 Micro-C 400M	<ul style="list-style-type: none"><li>• <a href="https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_400M_R1.fastq">https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_400M_R1.fastq</a></li><li>• <a href="https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_400M_R2.fastq">https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_400M_R2.fastq</a></li></ul>
GM12878 Micro-C 800M	<ul style="list-style-type: none"><li>• <a href="https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_800M_R1.fastq">https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_800M_R1.fastq</a></li><li>• <a href="https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_800M_R2.fastq">https://s3.amazonaws.com/dovetail.pub/HiC/fastqs/MicroC_800M_R2.fastq</a></li></ul>

## 1.10 Support

For help or questions related please open a new issue on the github repository or send an email to: [support@dovetail-genomics.com](mailto:support@dovetail-genomics.com)



## INDICES AND TABLES

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