

**MA5112**

**Group 1 Nextflow Script: Hi-C Assay**

Gavin Farrell

Sarah Drury

Aodán Laighneach

Eric Tornabell

&lt;&gt; Edit file Preview changes

Tabs 8 No wrap

```
1  #!/bin/bash/env nextflow
2
3  /*
4   * Written by Sarah
5   */
6
7  /*
8   Perform Hi-C analysis on raw fastq.gz files. Adapted from https://github.com/nf-core/hic/
9   Run script using "nextflow run hi-c.nf" Optionally, add "--restriction <sequence>" to
10  specify the restriction enzyme cut site. Default is 'A^AGCTT' (digestion by HindIII),
11  which was the protocol run on the test dataset provided.
12  */
13
14  // SETUP STEPS
15
16  // Declare pipeline parameters
17  reads = "$baseDir/data/*_R{1,2}.fastq.gz"
18  reference = "$baseDir/reference/*.fsa"
19  scripts = "$baseDir/scripts"
20  outDir = "$baseDir/results"
21  project_dir = projectDir
22
23  /*
24   Default restriction cutting site. Can be overridden by
25   using the --'<sequence>' command when launching the script
26  */
27
28  params.restriction = 'A^AGCTT'
29
30
31  // Set up channels for raw (.fastq.gz) and reference genome (.fsa) files
32  Channel
33    .fromPath(reads)
34    .map{file -> [file.simpleName, file]}
35    .into{ reads_ch; reads_for_fastqc }
36
37  Channel
38    .fromPath(reference)
39    .map{file -> [file.simpleName, file]}
40    .into{ reference_ch; fasta_for_resfrag_ch; fasta_for_chromsize }
41
42
43  // STEP 1: Align reads to genome
44
45  // Step by Bowtie2 to index the reference genome
46  process create_bowtie2_index {
47    publishDir path: "$outDir/index", mode: "copy"
48
49    input:
50      tuple val(base), file(ref) from reference_ch
51
52    output:
53      tuple val(base), file("*.bt2") into index_ch
54
55    script:
56      """
57      bowtie2-build ${ref} ${base}
58      """
59  }
60
61  /*
62   Process cuts the reference genome at the restriction enzyme cutting site. The digested reference genome is
63   used in the get_valid_interaction process to remove experimental artefacts.
64  */
65  process get_restriction_fragments {
66    publishDir path: "$outDir/index", mode: "copy"
67  }
```

```

68     input:
69     file(fasta) from fasta_for_resfrag_ch
70
71     output:
72     file("*.bed") into res_frag_ch
73
74     script:
75     """
76     python $scripts/digest_genome.py \
77         -r ${params.restriction} \
78         -o restriction_fragments.bed ${fasta[1]}
79     """
80 }
81
82
83 // R1 and R2 reads are aligned separately to the reference genome.
84 process bowtie2_end_to_end {
85     publishDir path: "$outDir/align", mode: "copy"
86
87     input:
88     each(reads) from reads_ch
89     tuple val(base), file(index) from index_ch
90
91     output:
92     tuple val(prefix), file("${name}.bam") into end_to_end_bam
93
94     script:
95     name = reads[0]
96     prefix = name.toString() - ~/(_R1|_R2|_val_1|_val_2|_1|_2)/
97
98     """
99     bowtie2 -x ${base} \
100         -U ${reads[1]} \
101         -S ${reads[0]}.bam \
102         --very-sensitive \
103         -L 30 \
104         --score-min L,-0.6,-0.2 \
105         --end-to-end \
106         --reorder
107     """
108 }
109
110 // R1 and R2 reads are combined to produce a paired-end bam.
111 process combine_mapped_files{
112     publishDir path: "$outDir/align", mode: "copy",
113     saveAs: {filename ->
114         filename.indexOf(".pairstat") > 0 ? "$outDir/stats/${filename}" : "${filename}"
115     }
116
117     input:
118     tuple val(sample), file(aligned_bam) from end_to_end_bam.groupTuple()
119
120     output:
121     tuple val(sample), file("${sample}_bwt2pairs.bam") into paired_bam
122     tuple val(online), file("*.pairstat") into all_pairstat
123
124     script:
125     online = sample.toString() - ~/(\.[0-9]+)/
126
127     """
128     python $scripts/mergeSAM.py \
129         -f ${aligned_bam[0]} \
130         -r ${aligned_bam[1]} \

```

---

```

129         -r ${aligned_bam[1]} \\\
130         -o ${sample}_bwt2pairs.bam \\\
131         --single --multi -t
132     """
133 }
134
135 /*
136  * Written by Gavin
137  */
138
139 // STEP 2: Detection of valid interaction products
140
141 // Sets up and defines the process for getting valid interaction (GVI). |
142 process get_valid_interaction{
143     publishDir path: "$outDir/valid", mode: "copy",
144     saveAs: {filename ->
145         filename.indexOf(".RSstat") > 0 ? "$outDir/stats/$filename" : "$filename"}
146
147     // Defines the input values for the GVI process
148     input:
149         tuple val(sample), file(pe_bam) from paired_bam
150         file(frag_file) from res_frag_ch.collect()
151
152     // Defines the output values for the GVI process
153     output:
154         tuple val(sample), file("*.validPairs") into valid_pairs
155         tuple val(sample), file("*.validPairs") into valid_pairs_4cool
156         tuple val(sample), file("*.RSstat") into all_rsstat
157
158     // Defines the python script for the GVI process to utilise
159     script:
160         """
161         python $scripts/mapped_2hic_fragments.py \\\
162             -f ${frag_file} \\\
163             -r ${pe_bam} \\\
164             sort -T /tmp/ -k2,2V -k3,3n -k5,5V -k6,6n
165         """
166 }
167
168 // STEP 3: Duplicates removal
169
170 // Sets up and defines the process for removing duplicate (RD) reads in the data
171 process remove_duplicates {
172     publishDir path: "$outDir/valid", mode: "copy",
173     saveAs: {filename ->
174         filename.indexOf(".mergestat") > 0 ? "$outDir/stats/$filename" : "$filename"}
175
176     // Defines the input values for the RD process
177     input:
178         tuple val(sample), file(vpairs) from valid_pairs.groupTuple()
179
180     // Defines the output values for the RD process
181     output:
182         tuple val(sample), file("*.allValidPairs") into all_valid_pairs
183         tuple val(sample), file("*.allValidPairs") into all_valid_pairs_4cool
184         file("*.") into all_mergestat
185
186     // Defines the bash script code used for the RD process to operate
187     script:
188         """
189         ## Sort valid pairs and remove read pairs with same starts (i.e duplicated read pairs)
190         sort -T /tmp/ -S 50% -k2,2V -k3,3n -k5,5V -k6,6n -m ${vpairs} | \
191         awk -F"\t" 'BEGIN{c1=0;c2=0;s1=0;s2=0}{c1!=\2 || c2!=\5 || s1!=\3 || s2!=\6}{print;c1=\2;c2=\5;s1=\3;s2=\6}' >> ${sample}.allValidPairs
192         echo -n "valid_interaction\t" >> ${sample}.allValidPairs.mergestat
193         cat ${vpairs} | wc -l >> ${sample}.allValidPairs.mergestat
194         echo -n "valid_interaction_rmdup\t" >> ${sample}.allValidPairs.mergestat
195         cat ${sample}.allValidPairs | wc -l >> ${sample}.allValidPairs.mergestat

```

```

196     ## Count short range (<20000) vs long range contacts
197     awk 'BEGIN{cis=0;trans=0;sr=0;lr=0} \2 == \5{cis=cis+1; d=\$6>\$3?\$6-\$3:\$3-\$6; if (d<=20000){sr=sr+1}else{lr=lr+1}} \2!=\5{trans=trans+1}END{print "trans_interact
198     ""
199 }
200
201
202 /*
203  * Written by Aodán
204  */
205
206 // STEP 4: Generate raw and normalized contact maps
207
208 /*
209  *Chromosome/scaffold sizes must be provided to build contact maps. First, samtools is used
210  *to index reference fasta to .fai file. Columns 1 & 2 (containing chromosome/scaffold) I.D. and
211  *length respectively will be cut to chrom.size file for future use.
212  */
213
214 process make_chrom_size {
215     publishDir path: "$outDir/chrom_size", mode: "copy"
216
217     input:
218     tuple val(base), file(fasta) from fasta_for_chromsize
219
220     output:
221     file("*.size") into chromosome_size, chromosome_size_cool
222
223     script:
224     """
225     samtools faidx ${fasta}
226     cut -f1,2 ${fasta}.fai > chrom.size
227     """
228 }
229
230 /*
231  *External program file build_matrix is used to produce contact map(s). This will require
232  *inputs of resolution, chromosome length and file denoting valid interaction pairs. Valid pairs
233  *and chromosome lengths will be collected from established channels above. Resolution will be
234  *provided through bin_size parameter, where two resolutions will be considered. Output will consist of
235  *genomic interval files (.BED) consistent with resolution and raw (unnormalised) matrix file.
236  */
237
238 // Resolutions for contact maps
239 bin_size = '1000000,500000'
240 map_res = bin_size.tokenize(',')
241
242 process build_contact_maps {
243     publishDir path: "$outDir/matrix/raw", mode: "copy"
244
245     input:
246     tuple val(sample), file(vpairs), val(mres) from all_valid_pairs.combine(map_res)
247     file chrsize from chromosome_size.collect()
248
249     output:
250     file("*.matrix") into raw_maps
251     file "*.bed"
252
253     script:
254     """
255     ${scripts}/build_matrix --matrix-format upper \\\
256         --binsize ${mres} \\\
257         --chrsize ${chrsize} \\\
258         --ifile ${vpairs} \\\
259         --oprefix ${sample}_${mres}
260     """
261 }
262
263 ~~~

```

```

262
263
264 /*
265  *Biases in the raw matrix file will be normalised using ICE. Raw matrix file will be used from
266  *previously-established raw_maps channel. Low and high counts will be prefiltered before matrix is
267  *normalised. Minimum iterations for normalisation will be set at 100. Output consists of normalised (iced)
268  *matrices for specified resolutions.
269  */
270
271 process run_ice{
272     publishDir "$outDir/matrix/iced", mode: "copy"
273
274     input:
275         file(rmaps) from raw_maps
276
277     output:
278         file("iced.matrix") into iced_maps
279
280     script:
281         prefix = rmaps.toString() - ~/(\.matrix)?$/
282         ""
283         ice --filter_low_counts_perc 0.02 \
284             --results_filename ${prefix}_iced.matrix \
285             --filter_high_counts_perc 0 \
286             --max_iter 100 \
287             --eps 0.1 \
288             --remove-all-zeros-loci \
289             --output-bias 1 \
290             --verbose 1 ${rmaps}
291         ""
292     }
293
294
295 /*
296  * Written by Eric
297  */
298
299 // STEP 5: Generate statistics files and quality control report
300 /*
301  *Combines the R1 and R2 statistics files. Statistics about read pairs filtering are available in the
302  *.mRSstat file, and pairing statistics are available in the .mpairstat file.
303  */
304
305 process merge_stats {
306     publishDir "$outDir/mstats", mode: "copy"
307
308     input:
309         tuple val(prefix), file(fstat) from all_rsstat.groupTuple().concat(all_pairstat.groupTuple())
310
311     output:
312         file("") into all_mstats
313
314     script:
315         sample = prefix.toString() - ~/(_R1_|_R2|_val_1|_val_2|_1|_2)/
316         if ( (fstat =~ /.pairstat/) ){ ext = "mpairstat" }
317         if ( (fstat =~ /.RSstat/) ){ ext = "mRSstat" }
318         ""
319         python $scripts/merge_statfiles.py -f ${fstat} > ${prefix}.${ext}
320         ""
321     }
322
323
324
325 // Perform quality control using FastQC
326 process fastqc {
327     input:
328         tuple val(sample), file(reads) from reads_for_fastqc
329
330

```

```

329
330     output:
331     file ".*_fastqc.{zip,html}" into fastqc_results
332
333     script:
334     """
335     fastqc -q ${reads}
336     """
337 }
338
339 /*
340 Finally MultiQC process is a reporting tool that parses summary statistics from results and log files generated by other
341 bioinformatics tools. It recursively searches through all provided file paths and finds files that it recognizes. It parses
342 relevant information from these and generates a single stand-alone HTML report file.
343 */
344 process multiqc {
345     publishDir "$outDir", mode: "copy"
346
347     input:
348     file (``) from fastqc_results.collect().ifEmpty([])
349
350     output:
351     file "multiqc_report.html" into multiqc_report
352     file "multiqc_data"
353
354     script:
355     """
356     multiqc .
357     """
358 }

```