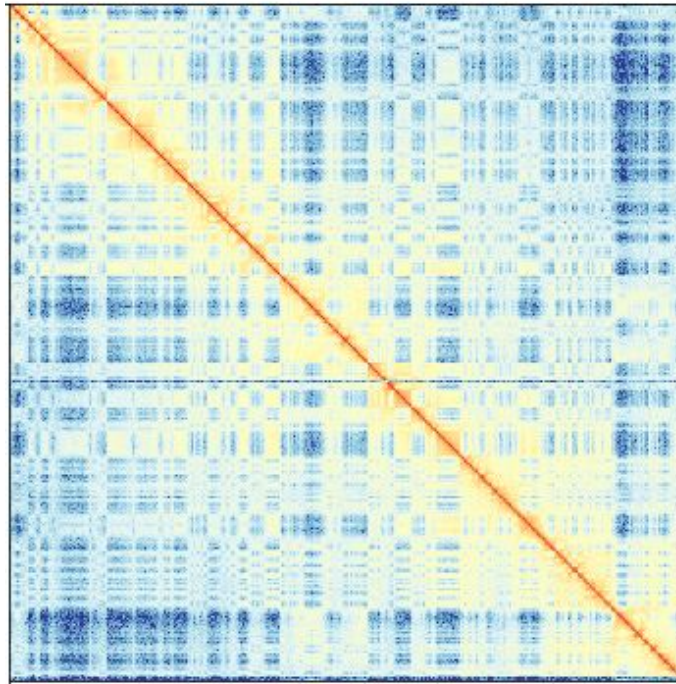


Projet PolledHiC

Stage Facultatif de Master 1 - Bioinformatique et Biologie des Systèmes



Thomas Faraut, Alain Pinton

- GenPhySE – INRAE Occitanie -

Jérémy Martin

01/07/2020 - 31/08/2020

- Master Bioinformatique et Biologie des Systèmes -

- Université Paul Sabatier (Toulouse) -



INRAE



Je tiens à remercier mon responsable pédagogique, Roland Barriot, qui a su se montrer très disponible en amont de ce stage pour régler le flot continu de problèmes administratifs malgré des conditions de travail compromises.

Je tiens également à remercier mes encadrants, Thomas Faraut et Alain Pinton qui ont défendu le projet pour que le stage soit maintenu, ainsi que pour leur disponibilité, leur écoute, et leurs encouragements.

Enfin je remercie toutes les personnes qui ont eu la patience d'écouter mes "histoires de vaches sans cornes" quand ces dernières occupaient toutes mes pensées.

Jérémy Martin

PolledHiC

PolledHiC

Introduction

La méthode HiC :

Les données SeqOccin :

NF-Core/HiC

Workflow

MultiQC

Tableau récapitulatif des données :

Les matrices de contacts

Conda : polledHiC_env

hicExplorer

ConvertFormat

SumMatrices

AdjustMatrix

Normalize et CorrectMatrix

Diagnostic plot

Carte HiC

Recherche de TADs

hicFindTADs :

hicPlotTADs :

Snakemake et automatisation

nfcore-hic

nfcore-hic.config

nfcore-hic.sh

hicexplorer

config.yaml

make_sampled_protocol.py

Snakefile

hicexplorer.sh

TL;DR : Comment lancer une analyse sur Genologin :

Introduction

Le projet PolledHiC est un projet réalisé dans le cadre d'un stage facultatif de Master 1 Bioinformatique et Biologie des systèmes.

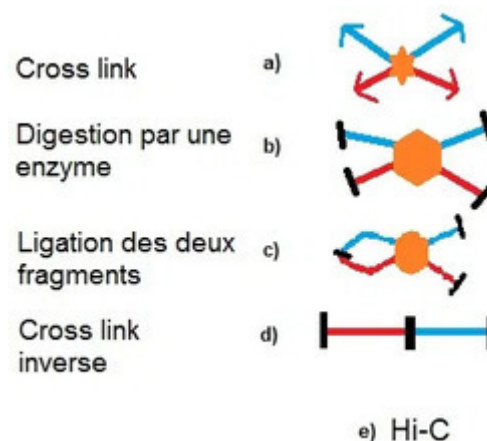
Il a été initié par Thomas FARAUT et Alain PINTON (INRAE Occitanie - UMR GenPhySE).

Le projet consiste à traiter les données de séquençage HiC de l'équipe de SeqOccin pour tenter d'obtenir des cartes chromosomiques de différentes régions d'intérêts chez les bovins. Le but est de fournir un pipeline permettant l'analyse et la comparaison des structures tri-dimensionnelle entre les phénotypes 'avec cornes' et 'sans cornes' chez certains d'entre eux.

La méthode HiC :

La méthode de séquençage HiC (High Chromosome Contact map), est une technique de séquençage haut débits permettant d'étudier les interactions au niveau génétique.

Les étapes sont les suivantes :



- Cross-link pour 'figer' les chromosomes.
- Digestion par des enzymes de restrictions.
- Ligation des fragments cross-linker.
- Élimination du 'cross-link'
- Filtrage des séquences sur les régions d'intérêt
- Visualisation (cartes HiC et plot TADs)

Les données SeqOccin :

Les données ont été obtenus par le biais du projet « **Séquençage Occitanie Innovation** » (soit SeqOccin).

Elles contiennent les informations de deux trio de bovins pour lesquels ont été réalisés :

- différents protocoles HiC (Arima, Dovetail, PhaseG et Maison).
- différents runs avec différent niveaux de coverage.

Dans le cadre de ce projet, les protocoles PhaseG n'ont pas été retenus.

NF-Core/HiC

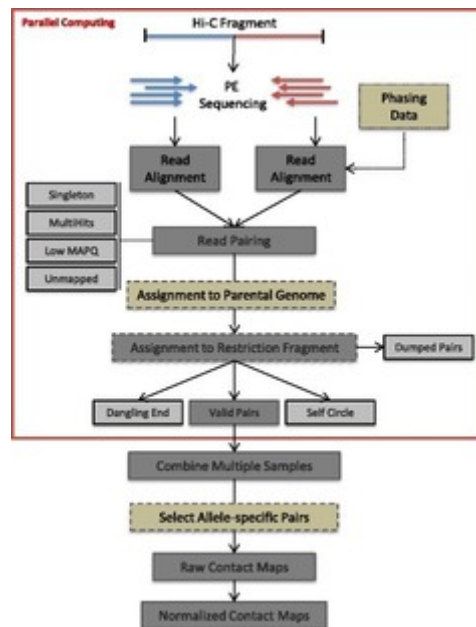
Le pipeline [NF-Core/HiC](#) est basé sur le workflow HiC-Pro. Il permet d'obtenir des matrices de contacts à partir de données brutes au format FastQ.

A l'issue de ces étapes, les manipulations sur les matrices de contacts seront réalisés via le logiciel [HiCExplorer](#) .

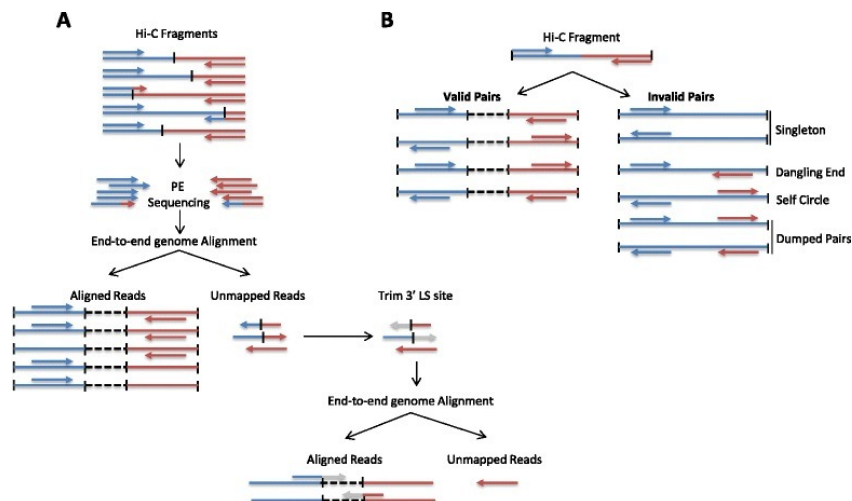
Workflow

Le workflow de NF-Core/hic se réalise en X étapes :

1. Alignements des reads
2. Appariements des reads
3. **Mapping**
4. Obtention de la matrice de contacts brut



Au cours de ce pipeline, l'étape de mapping va déterminer les "**valid-pairs**". C'est à partir de ces valid-pairs que la matrice de contact pourra être créée, l'étape de mapping est donc très importante.

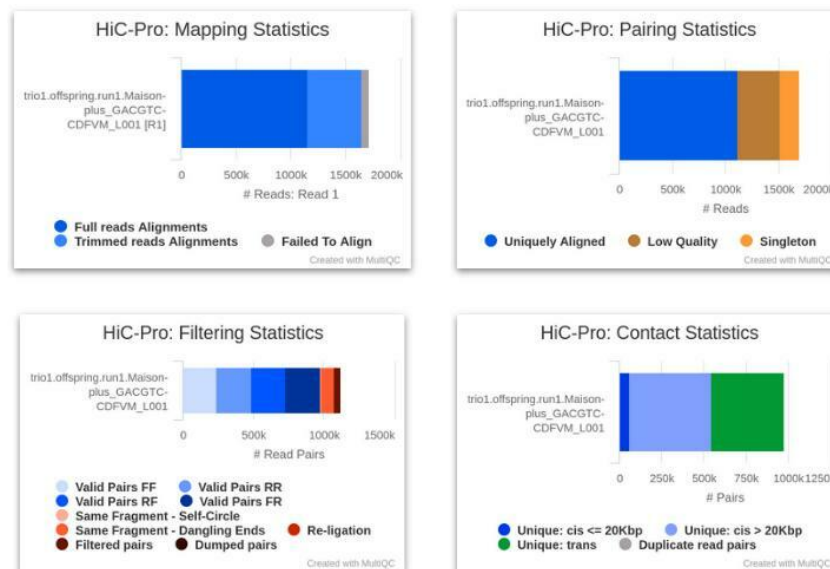


Tout ce qui n'est pas considéré comme une valid-pairs n'est pas retenu dans l'élaboration de la carte de contact. Ils peuvent être obtenus en première intention à l'issue des premières étapes d'alignements, ou en seconde intention, après une étape de "trim" des fragments de restrictions / ligations.

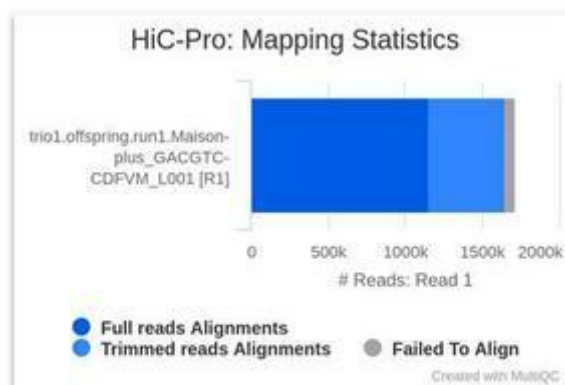
MultiQC

À l'issue de son processus, NF-Core/hic permet également d'obtenir un rapport [MultiQC](#).

MultiQC est un outil permettant de présenter sous forme de rapport html les résultats de nombreux outils bio-informatiques. Ici ces rapports contiennent les résultats statistiques des étape de mapping et de construction de la carte de contact.



MAPPING STATISTICS :



Full reads alignment = 67.2% ; Trimmed reads alignment = 28.9% ;
Failed to align = 3.9%

Il est attendu qu'une grosse majorité des reads soit directement aligné sur le génome (80 à 90%). Les 'trimmed reads alignments' correspondent aux reads alignées en seconde intention. Il s'agit de fragments chimériques qui sont détectés avec le fragments de ligation. Un niveau trop élevée de ces fragments chimériques peut être le reflet d'un problème lors de l'étape de ligation.

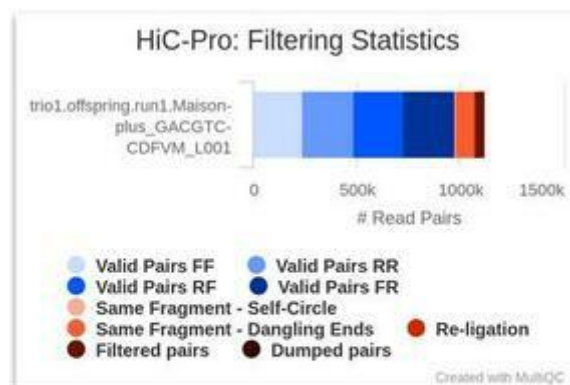
PAIRING STATISTICS :



Uniquely aligned = 65.8% ; Low quality = 23.3% ;
Singleton = 10.9%

Les reads sont censé être "paired". Le taux des fractions dépend grandement de la complexité du génome et de la fraction de reads non mappée. On s'attend à ce que la fraction de singleton soit proche de la fraction de reads non mappé. Un taux de singleton trop élevée et/ou différent du nombre de reads non mappé peut être révélateur d'un problème avec les données.

FILTERING STATISTICS :

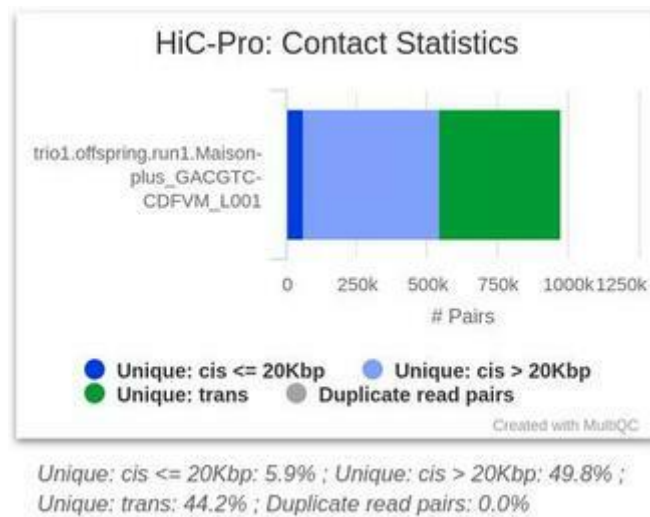


Valid Pairs FF: 21.8% ; Valid Pairs RR: 21.8% ;
Valid Pairs RF: 21.5% ; Valid Pairs FR: 22.3% ;
Same Fragment - Same Circle: 0.3% ;
Same Fragment - Dangling Ends: 8.1% ;
Re-ligation: 0.8% ; Filtered pairs: 3.5% ; Dumped pairs: 0.0%

La qualité de l'analyse peut être évaluée en filtrant les paires de reads.

Comme la ligation est un processus aléatoire, on attend une répartition de 25% pour chaque classe de "**valid-pairs**" (*Forward-Forward, Forward-Reverse, Reverse-Forward* et *Reverse-Reverse*). De la même façon, une fraction élevée de *dangling-end* ou de *self-circle read pairs* est représentatif d'une faible qualité expérimentale et peut révéler un problème lors des étapes de digestion, d'élongation ou de ligation.

CONTACT STATISTICS :



Il est généralement attendu d'observer une distribution centrée autour de la taille de l'insert utilisé. Il est normal d'observer une majorité de contact type cis- éloigné (plus de 20kb de distance sur un même chromosome) ou trans- (deux régions sur deux chromosomes différents).

La fraction des doublons est également présentée. Un niveau élevé de duplication indique une faible complexité moléculaire et un biais de PCR potentiel.

Tableau récapitulatif des données :

Voici ci dessous, un tableau récapitulatif des données utilisées :

trio	animal	run	protocol	paired reads	valid pairs
trio1	father	run1	Arima	89637817	22279893
trio1	father	run1	Dovetail	3700008	1866639
trio1	father	run1	Maison	141875156	66790091
trio1	father	run2	Arima	119832383	30993207
trio1	mother	run1	Arima	76366518	16817078
trio1	mother	run1	Dovetail	108323412	49112300
trio1	mother	run1	Maison	134238122	82351532
trio1	mother	run2	Arima	105484493	0
trio1	offspring	run1	Arima	1688139	387688
trio1	offspring	run1	Dovetail	1496149	437911
trio1	offspring	run1	Maison-minus	1541573	613689
trio1	offspring	run1	Maison-plus	1715720	976474
trio1	offspring	run2	Arima	130341495	30135249
trio1	offspring	run2	Dovetail	164244926	49182828
trio1	offspring	run2	Maison-plus	67743665	35536687
trio1	offspring	run3	Maison-plus	188344770	110398242
trio2	father	run1	Maison	128040893	11267457
trio2	father	run2	Maison	100015296	8456432
trio2	mother	run1	Maison	159863877	80781661
trio2	mother	run2	Maison	113850699	54282196
trio2	offspring	run1	Maison	141212982	75281563
trio2	offspring	run2	Maison	107893352	54602426

Au vu des informations précédentes, il a été établi que les paired reads et les valid-pairs sont déterminant dans la qualité de l'analyse qui en découle.

Il est observable que le trio1 mother run2 Arima n'a aucun valid-pairs. Il a été décidé que les deux runs ne seraient pas pris en compte. Il y avait donc 20 runs qui composait le matériel d'étude.

Il est également observable que les individus du trio2 présentent moins de données que le trio1. Cela se ressentira lors de l'analyse.

Les matrices de contacts

Les matrices de contacts se présentent sous la forme :

```
1  A  B  10
2  A  C  23
3  B  C  24
4  (...)
```

Ce format permet d'optimiser l'espace mémoire en ne tenant pas compte des interactions nulle (aucune valeur à 0). Il est également compatible avec d'autre logiciel d'analyse tel que HiTC Bioconductor package ou HiCPlotter software.

Cependant, ce format n'est pas directement pris en charge par hicExplorer, mais il est possible de le convertir au format hdf5 reconnu par l'outil.

A ce stade, il y a une matrice de contact par run et par résolutions étudiées (200000, 50000, 25000 et 10000 bins). Dans le cadre de l'étude, il y a donc 80 matrices de contacts.

Conda : polledHiC_env

Pour réaliser le reste des analyses, il est nécessaire d'utiliser la dernière version d'hicExplorer.

Pour s'assurer de la répétabilité de la manipulation, il est conseillé de travailler sous un environnement conda identique à celui de l'étude :

```
1  conda env create -n polledHiC_env python=3.7.8
2  conda activate polledHiC_env
3
4  conda install hicexplorer=3.5.1 snakemake=5.14.0
```

Une autre façon est de télécharger [ce fichier](#)
`polledHiC_env.txt` , puis :

```
1 conda create --name polledHiC_env --file polledHiC_env.txt
2
3 conda activate polledHiC_env
```

Dans les deux cas, il est indispensable d'activer l'environnement de travail pour que les étapes suivantes soient correctement exécuter.

hicExplorer

hicExplorer est une suite d'outils permettant l'analyse de données HiC.

Le pipeline d'analyse réalisé utilise de nombreux outils et permet d'obtenir, à l'issue de ce dernier, une carte chromosomique à différente échelle centrée sur les 3 premiers millions de bases du chromosome 1.

ConvertFormat

La première étape consiste à convertir les matrices hicpro vers un format hdf5 exploitable par le reste des outils d'hicExplorer :

```
1 hicConvertFormat --matrices {input.raw} --bedFileHicpro
  {input.bed} --outFileName {output} --inputFormat hicpro --
  outputFormat h5
```

Le format HDF5 est un conteneur de fichier construit en ensemble de données. Cela donne aux fichiers HDF5 une structure en arborescence :

```
— matrix [HDF5 group]
  |— barcodes
  |— data
  |— indices
  |— indptr
  |— shape
  |— features [HDF5 group]
    |— _all_tag_keys
    |— feature_type
    |— genome
    |— id
    |— name
    |— pattern [Feature Barcoding only]
    |— read [Feature Barcoding only]
    |— sequence [Feature Barcoding only]
```

SumMatrices

La seconde étape consiste à sommer les matrices de même résolutions pour obtenir un maximum d'informations pour chaque individus.

Cela se fait en deux temps :

```
1 # Merge runs by protocols
2 hicSumMatrices --matrices {input.indiv.protocol.h5} --outFileName
  {output.indiv.merged.protocol.h5}
3
4 # Merge protocols by indiv
5 hicSumMatrices --matrices {input.indiv.merged.protocol.h5} --
  outFileName {output.indiv.merged.h5}
```

AdjustMatrix

La commande hicAdjustMatrix permet de s'affranchir des contigs dans les données à analyser, et de se concentrer uniquement sur les chromosomes complets.

```
1 hicAdjustMatrix --matrix {input.indiv.merged.h5} --outFileName
  {output.indiv.adjusted.h5} --chromosomes 1 2 3 4 5 6 7 8 9 10 11
  12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 X --action
  keep
```

Normalize et CorrectMatrix

Les développeurs d'hicExplorer recommandent de réaliser une normalisation avant la correction de la matrice (*cf. la documentation de [hicNormalize](#) sur [readthedocs.io](#)*).

La normalisation des matrices de données permet de normaliser au plus petit nombre de lecture donné de toutes les matrices ou à une plage de 0 à 1.

```
1 hicNormalize --matrices {input.indiv.adjusted.h5} --outFileName
  {output.indiv.normalized.h5} --normalize smallest
```

La correction quant à elle permet d'éliminer les biais induit par les régions riches en GC, et normaliser le nombre de sites de restriction par bins. Étant donné qu'une fraction des bins peuvent être issue de régions fortement répétées, ces bins contiennent peu de contact et il est nécessaire de les filtrer.

```
1 hicCorrectMatrix correct --matrix {input.indiv.normalized.h5} --  
  outFile {output.indiv.corrected.h5}
```

A l'issue de cette étape, il a été obtenu pour chaque individus, quatre matrices finales (une pour chaque résolutions) prêtes à être analyser.

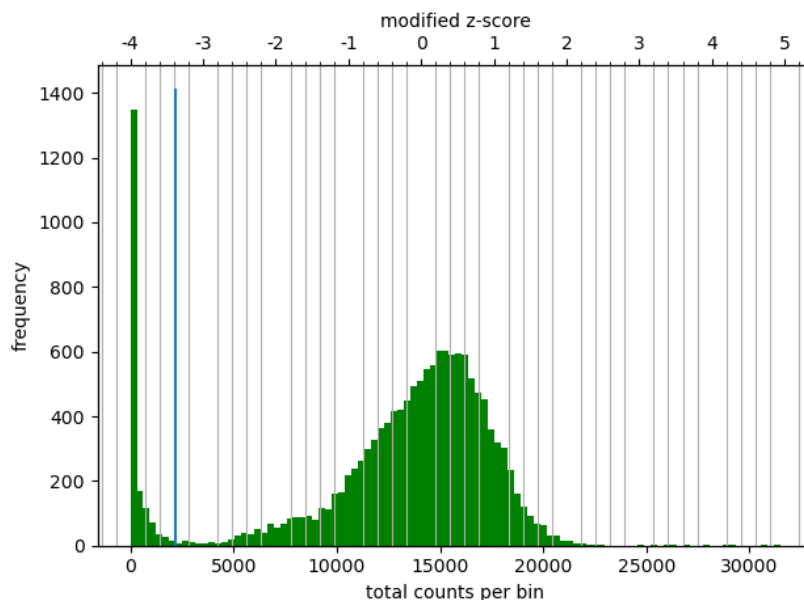
Diagnostic plot

Le diagnostic plot est un histogramme de la fréquence de comptage par bins. Cet histogramme permet de comparer les différentes résolutions sur un même jeu de données. Il est attendu d'observer une répartition Normale.

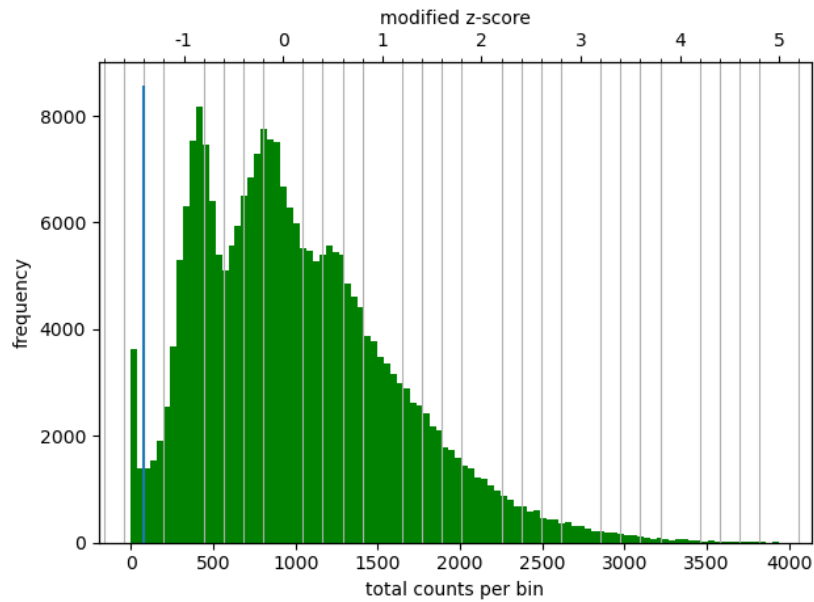
Le modified Z-Score est un indicateur permettant d'identifier les outliers de nos données. Le trait bleu sur ces histogrammes marque le MAD Z-Score : Median Absolute Deviation. C'est la valeur par défaut prise lors de la correction pour éliminer ces outliers.

Le choix de la bonne résolution est cruciale dans la qualité d'analyse, et dépend des données analysées :

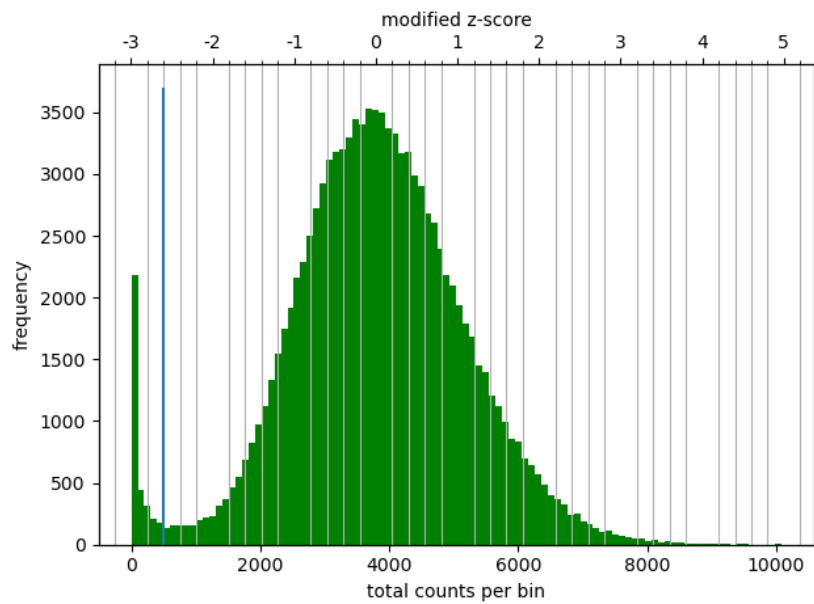
Avec une résolution à 200.000 pb par bins, on observe chez Trio1 Father une répartition qui est loin de s'apparenter à une Gaussienne.



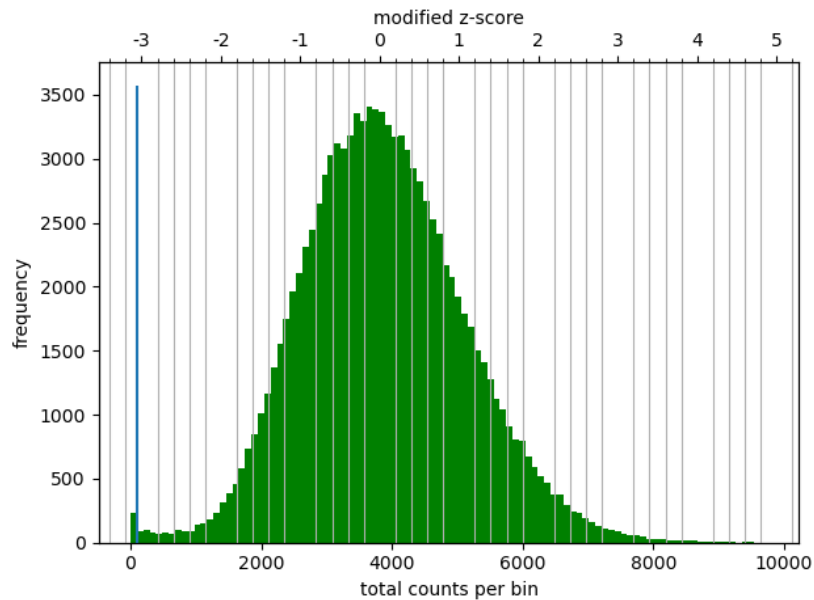
De la même manière, avec une résolution trop fine, certains jeu de données peuvent donner des résultats incompatibles avec la suite des analyses. Ici on observe l'individu Trio2 Offspring à une résolution de 10.000 pb par bins :



En choisissant une résolution adéquate avec les données, on retrouve alors une répartition s'apparentant à une Normale, comme par exemple ici, avec Trio1 Offspring à 25k bins :



On observe un gros pics à 0. Les étapes d'ajustement et de correction permettent d'éliminer ces biais. Voici le diagnostic plot du même individu, à la même résolution, après réalisation de ces étapes :

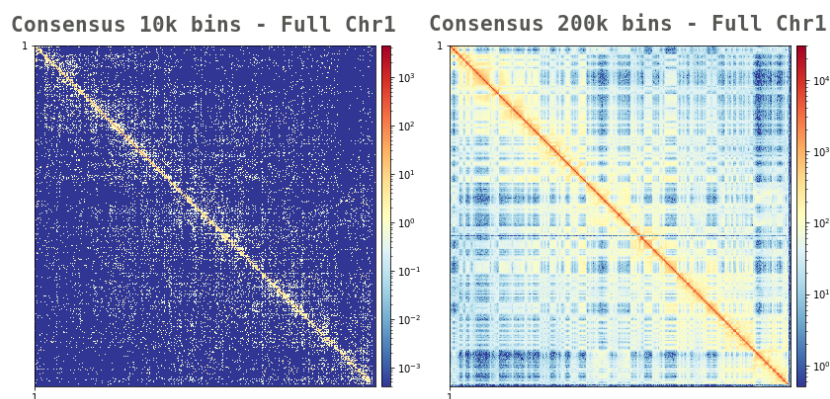


Carte HiC

Enfin, il devient alors possible de créer les cartes de contacts des régions qui nous intéressent à partir des résultats précédents.

Dans le but d'établir un profil, il a été créé un individus **"Consensus"** issue de la somme des matrices finales de chaque individus.

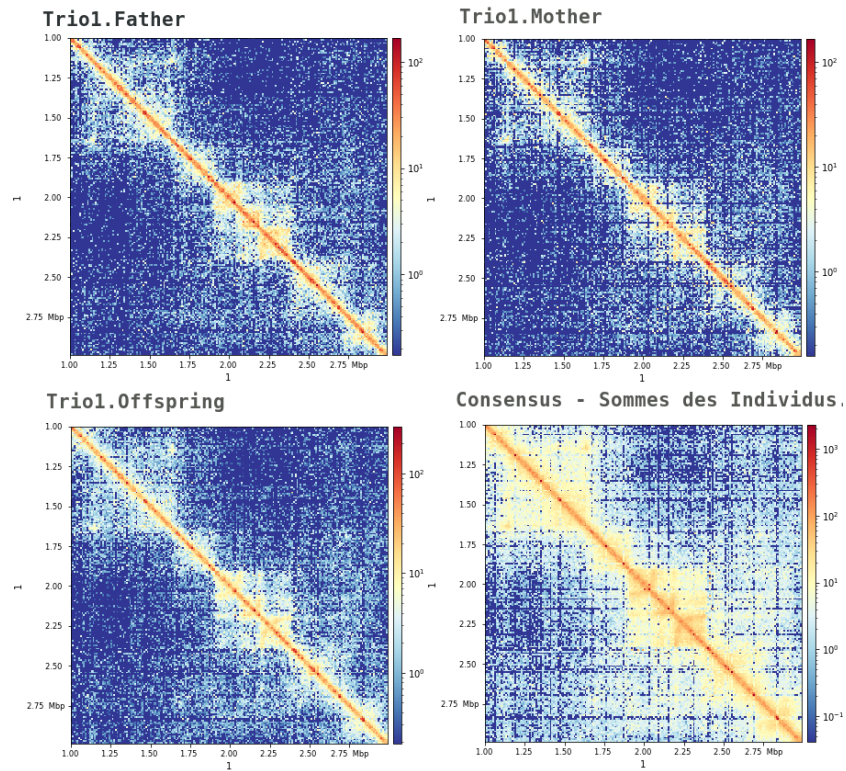
La résolutions d'études va également dépendre de la fenêtre qui va nous intéresser. Plus cette fenêtre est grande, moins on doit être résolutif :



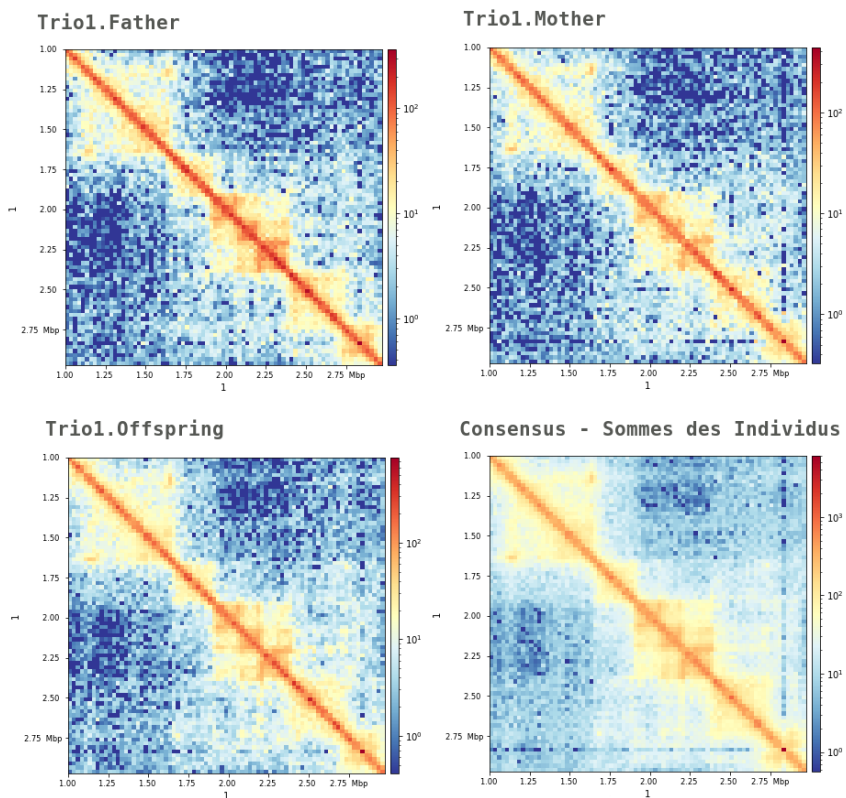
La fenêtre choisie va des positions 1M à 3M de bases sur le chromosome 1. Nous nous concentrerons donc sur les plus fortes résolutions : 10.000 pb et 25.000 pb par bins.

Cette fenêtre a été choisie pour prendre en compte les mutations Celtic, Friesian, Mongolian et Guarani situé entre les positions 2M4 et 2M7.

10k bins resolutions :



25k bins resolutions :



Au vu des données possédés, il semblerait que la résolution 25k bins soit la plus adéquate pour analyser nos données. On notera toute fois que la mise en place de l'individu "*Consensus*" a permis d'obtenir des résultats plus probant à la résolution 10k bins que les individus.

On observe alors la formation de **TADs** (*Topologically Associated Domains*) sur ces cartes. Ces TADs sont représentés par des "carrés" sur la diagonale de la carte, notamment entre les positions 1.20 et 1.65 , 1.70 et 1.95 , 2.00 et 2.40 et enfin entre 2.40 et 2.75.

Il semblerait également que la matrices consensus permette de retrouver ces patterns à la résolution 10k bins contrairement aux matrices issue des individus.

Cet interprétation n'étant qu'une interprétation graphique, il s'en suit donc une dernière étape : la recherche et la caractérisation des TADs.

Recherche de TADs

La recherche de TADs se fait via les outils `hicFindTADs` et `hicPlotTADs`

hicFindTADs :

La fonction `hicFindTADs` mesure le TAD-separation score pour identifier des degrés de séparation sur chaque bins de la matrice hic.

```
1 hicFindTADs --matrix {input} --outPrefix {outprefix} \  
2 --correctForMultipleTesting fdr --chromosomes 1
```

En sortie de cette fonction, il y a plusieurs fichiers produits :

- **{outprefix}_boundaries.bed** et **{outprefix}_boundaries.gff** : contiennent les informations sur les frontières de TADs.
- **{outprefix}_domains.bed** : contient les positions des TADs.
- **{outprefix}_score.bedgraph** et **{outprefix}_score.bm** : contiennent les TAD-separation score. Le fichier `score.bm` est une 'matrice bedgraph' qui peut être utiliser pour la représentation graphique des TAD-separation score dans `hicPlotTADs`.
- **{outprefix}_zscore_matrix.h5** : contient les matrices où les score ont été remplacés par les z-score (nécessaire pour la création des fichiers boundaries).

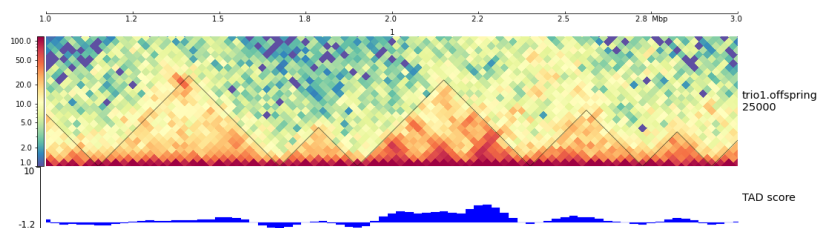
hicPlotTADs :

En plus des sortie de `hicFindTADs` , il est nécessaire de mettre en place un fichier `tracks.ini` pour plotter correctement les TADs.

Ces fichiers sont composé comme ci-dessous :

```
1  [x-axis]
2  fontsize=10
3
4  [hic matrix]
5  file = # path/to/matrix/corrected.h5
6  title = # title i.e. 'sample_resolution'
7  # depth is the maximum distance plotted in bp. In Hi-C tracks
8  # the height of the track is calculated based on the depth such
9  # that the matrix does not look deformed
10 colormap = Spectral_r
11 depth = 1500000
12 transform = log1p
13 file_type = hic_matrix
14
15 [tads]
16 file = # path/to/specific/{outprefix}_domains.bed from
17     hicFindTADs of hic matrix file.
18 file_type = domains
19 border_color = black
20 color = none
21 # the tads are overlay over the hic-matrix
22 # the share-y options sets the y-axis to be shared
23 # between the Hi-C matrix and the TADs.
24 overlay_previous = share-y
25
26 [bedgraph]
27 file = # path/to/specific/{outprefix}_score.bedgraph from
28     hicFindTADs of hic matrix file.
29 color = blue
30 height = 3
31 title = TAD score
32 max_value = 10
33
34 [spacer]
35 height = 1
```

Il est donc nécessaire de réaliser un fichier `tracks.ini` par plot que l'on souhaite réaliser.



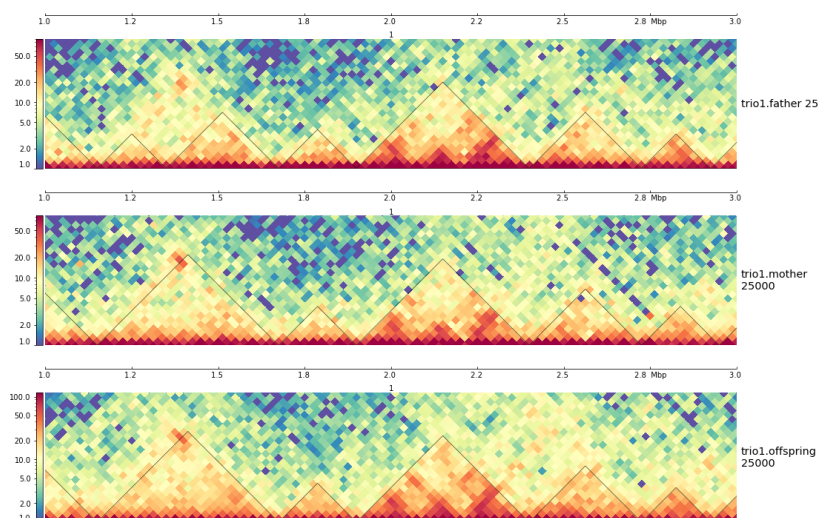
Il est aussi possible de mettre les plots directement les uns en dessous des autres par exemple :

```

1  [x-axis]
2  fontsize=10
3
4  [hic matrix]
5  file = hdf5/file1_corrected.h5
6  title = File1 25000 bins
7  # depth is the maximum distance plotted in bp. In Hi-C tracks
8  # the height of the track is calculated based on the depth such
9  # that the matrix does not look deformed
10 colormap = Spectral_r
11 depth = 1500000
12 transform = log1p
13 file_type = hic_matrix
14
15 [tads]
16 file = TADs/file1_domains.bed
17 file_type = domains
18 border_color = black
19 color = none
20 # the tads are overlay over the hic-matrix
21 # the share-y options sets the y-axis to be shared
22 # between the Hi-C matrix and the TADs.
23 overlay_previous = share-y
24
25 [spacer]
26 height = 1
27
28 [hic matrix]
29 file = hdf5/file2_corrected.h5
30 title = File2 25000 bins
31 colormap = Spectral_r
32 depth = 1500000
33 transform = log1p
34 file_type = hic_matrix
35
36 [tads]
37 file = TADs/file2_domains.bed
38 file_type = domains
39 border_color = black
40 color = none
41 overlay_previous = share-y
42
43 ...

```

Ici par exemple, les résultats du trio1 à la résolution 25.000 bins ont été assemblés dans le même fichier .jpg :



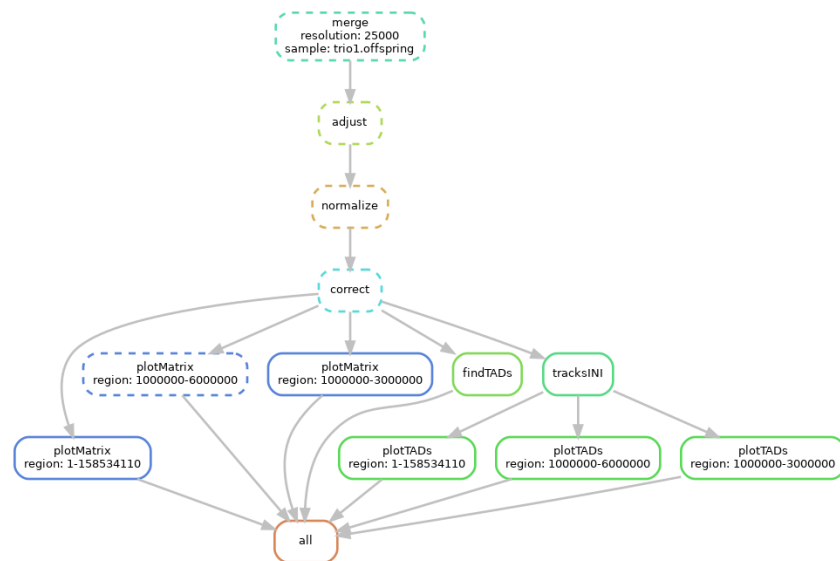
Il est alors possible d'observer que le TAD entre 1.15 et 1.75 chez mother et offspring se trouve découper en deux chez le père. Une approche génomique pourrait éventuellement montrer un événement d'indel à ce niveau expliquant la différence de TAD-separation score.

Snakemake et automatisation

Une fois la partie technique réalisée, il a fallu concentrer les efforts sur la mise en place d'une pipeline permettant la répétabilité, la reproductibilité, la portabilité et l'automatisation des analyses effectuées.

Snakemake est un gestionnaire de pipeline, issue de python et de GNU Make. Il est basé sur un système de "règles" (rules) qui possède un ou plusieurs input, et qui retourne au moins un output. Cet output peut être le produit final, comme un intermédiaire qui devient un input dans la règle suivante.

Par exemple, voici le DAG des étapes hicExplorer sur trio1.offspring, à la résolution 25.000 bins pour 3 régions d'intérêt :



Il est donc possible de décomposer l'intégralité des travaux réalisés en deux grandes étapes :

- nfcore-hic : pour traiter les données de séquençage et obtenir les matrices de contacts.
- hicexplorer : pour traiter les matrices de contacts et obtenir les cartes HiC et plotTADs

nfc core-hic

Après plusieurs tentatives pour intégrer le nextflow dans le Snakemake, il a été convenu de garder séparé les deux workflows. Le dossier nfc core-hic contient deux fichiers : **nfc core-hic.config** et **nfc core-hic.sh**.

nfc core-hic.config

Le fichier `nfc core-hic.config` contient les règles d'utilisation en temps et mémoire de certains sous-jobs pour la soumission sur le cluster. **Il n'est pas nécessaire de le modifier.**

nfc core-hic.sh

Le fichier `nfc core-hic.sh` est le script permettant d'initialiser l'espace de travail. Ce dernier **ne soumet pas le job sur le cluster**.

Il est nécessaire de personnaliser ce fichier, notamment les parties **Set parametres**, **Set restriction/ligation parametres** et **--bin-size** :

```
1  # Set parametres :
2
3  datadir= # path/to/reads_directories i.e. :
4           /work2/genphyse/dynagen/jmartin/polledHiC/data/reads_nfc core
5
6  wdir= # path/to/working_directory i.e. :
7        /work2/genphyse/dynagen/jmartin/polledHiC/workflows/nfc corehic_jm
8
9  genome= # path/to/genomic_bank i.e. :
10         /bank/bowtie2db/ensembl_bos_taurus_genome
11
12  chromsize= # path/to/chrom.length file i.e. :
13            /work2/genphyse/dynagen/jmartin/polledHiC/data/genome/chrom.len
```

Lors de l'utilisation du pipeline, l'architecture du répertoire data était la suivante

```
1  | polledHiC
2  | | data
3  | | | reads_nfc core
4  | | | | Arima
5  | | | | | trio1.mother.run1.Arima...fastq.gz
6  | | | | | trio1.mother.run2.Arima...fastq.gz
7  | | | | | trio1.offspring.run1.Arima...fastq.gz
8  | | | | | ...
9  | | | | Dovetail
10 | | | | | ...
11 | | | | ...
12 | | | genome
```

Dans la seconde partie, il est nécessaire de modifier le case [protocol] et renseigner les fragments de restriction et ligations :

```
1  # Set restriction/ligation parametres
2  case $prot in
3    *[protocol]) restriction= # restriction fragment as following
      : '^GATC,G^ANTC'
4    ligation= # ligation fragment as following :
      'GATCGATC,GANTGATC,GANTANTC,GATCANTC'
5    echo "run [protocol] : $run
6    res : $restriction
7    lig : $ligation";;
8    *) restriction='not recognized'
9    ligation='not recognized'
10   echo "run NOT RECOGNIZED : $run
11   res : $restriction
12   lig : $ligation";;
13  esac
```

Si il n'y a qu'un seul protocole à étudier, il est alors possible de remplacer le case par :

```
1  # Set restriction/ligation paramtres
2  restriction= # restriction fragment as following: '^GATC,G^ANTC'
3  ligation= # ligation fragment as following:
4  'GATCGATC,GANTGATC,GANTANTC,GATCANTC'
5  echo "run [protocol] : $run
6  res : $restriction
7  lig : $ligation"
```

Enfin, la ligne --bin-size contient les différentes résolutions d'études. Il est important de respecter la typologie d'écriture :

```
1  65 # most are default options - skipping ICE because normalizing
2  66 # each technical replicate is not relevant
3  67 nextflow run nf-core/hic \
4  68 -revision 1.1.0 \
5  69 -profile genotoul \
6  70 -name '$runid' \
7  71 -c '$conf' \
8  [...]
9  77 --reads '$readsdir'/*_R{1,2}.fastq.gz \
10 [...]
11 81 --restriction_site '$restriction' \
12 82 --ligation_site '$ligation' \
13 [...]
14 88 --bin_size # replace x,y,z as following : 'x,y,z' \
15 [...]
16 94 --skipCool' > $script
```

Par exemple, dans le cadre de l'étude réalisé, cette ligne contenais :

```
1  88 --bin-size '200000,50000,25000,10000'
```

Une fois ce script exécuté, dans chaque sous répertoires se trouvera un script permettant de soumettre le job sur le cluster :


```

1  |— polledHiC
2  |   |— nfcore-hic
3  |       |— Arima
4  |           |— Arima.sh
5  |       |— Dovetail
6  |           |— Dovetail.sh
7  |       |— Maison
8  |           |— Maison.sh
9  |       |— nfcore-hic.config
10 |       |— nfcore-hic.sh
11 |
12 |
13 |
14 # Exemple de soumission :
15 sbatch nfcore-hic/Arima/Arima.sh

```

hicexplorer

Pour pouvoir exécuter le pipeline, il est essentiel d'avoir activé l'environnement **polledHiC_env**.

Le répertoire hicexplorer contient quatre fichiers :

- config.yaml
- make_sampled_protocol.py
- Snakefile
- hicexplorer.sh

config.yaml

Fichier de configuration du Snakefile. Il est nécessaire de le personnaliser quasi-entièrement pour correspondre aux analyses souhaitées.

```

1  ### Config.yaml to set wildcards to Snakefile ###
2
3  workdir: # path/to/working_directory
4
5  sampled_protocols: sampled_protocols.tsv # don't edit this line
6  !
7
8  resolutions: # resolution of interest
9
10 chromosomes: # full chromosomes to adjust without contig
11
12 regions: # regions of interest
13
14 samples: # samples who matched with sample column of
15 sampled_protocols

```


i.e. : le config.yaml utilisé pendant ce projet

```
1  ### Config.yaml to set wildcards to Snakefile ###
2
3  workdir: /home/user/pollHiC/workflows/hicexplorer
4
5  sampled_protocols: sampled_protocols.tsv
6
7  resolutions: [200000,50000,25000,10000]
8
9  chromosomes: [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15,
10 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, X]
11
12 regions: ['1-158534110', '1000000-6000000', '1000000-3000000',
13 '1500000-3500000', '2000000-3000000', '2200000-2800000']
14
15 samples: [trio1.offspring, trio1.father, trio1.mother,
16 trio2.offspring, trio2.mother, trio2.father]
```

make_sampled_protocol.py

Ce script permet d'initier la création d'un fichier tabulé contenant les informations des différentes matrices de contact à traiter.

Il est nécessaire de renseigner la ligne 19 :

```
1  19 nfcore_dir = # 'path/to/nfcore_directory/*/' <- the '/*/' at
the end allow to only return directory with glob.
```

i.e. : le script utilisé pour l'analyse :

```
1  19 nfcore_dir = '/home/user/pollHiC/workflows/nfcore-hic/*/'
```

Snakefile

Le Snakefile est le fichier contenant toutes les règles d'exécution du snakemake. Ici non plus il n'y a pas de modification à apporter.

Le Snakefile contient 12 règles (ou 'rules') pour mener de bout en bout le pipeline d'analyse.

Il est appelée par le script hicexplorer.sh pour être lancer comme un job sur le cluster genologin. Il est toutefois possible de lancer un dry-run sur une session dédiée :

```

1  # Demande d'une session dédiée :
2  srun -c 1 --mem=8G --pty bash
3
4  # Activation de l'environnement conda si ce n'est pas déjà fait
5  :
6  conda activate polledHiC
7
8  # Dry-run :
9  snakemake -j 1 -p -n
10
11 # Plot du DAG en prenant en compte les informations du
12 config.yaml :
13 snakemake --dag | dot -Tpng -o dag.png

```

hicexplorer.sh

Le script hicexplorer.sh permet de soumettre l'analyse sur le cluster.

Par défaut, le Snakefile est paramétré pour réaliser un `dry-run`. Ce dernier liste les différentes étapes que va réaliser snakemake, sans les réaliser.

Cela permet d'avoir une première vision et de s'assurer que tout les fichiers de configurations soient correctement renseignés.

Si le dry-run s'exécute correctement, il suffit d'enlever l'option `-n` à la commande snakemake

```

1  # from
2  33 snakemake --jobs 1 -p -n\
3  34 --snakefile
4  /work2/genphyse/dynagen/jmartin/polledHiC/workflows/hicexplorer_j
5  m/Snakefile \
6  35 --directory
7  /work2/genphyse/dynagen/jmartin/polledHiC/workflows/hicexplorer_j
8  m/
9
10 # to
11 33 snakemake --jobs 1 -p \
12 34 --snakefile
13 /work2/genphyse/dynagen/jmartin/polledHiC/workflows/hicexplorer_j
14 m/Snakefile \
15 35 --directory
16 /work2/genphyse/dynagen/jmartin/polledHiC/workflows/hicexplorer_j
17 m/

```

Libre à vous de vouloir ou non laisser s'exécuter la création d'un **"consensus"** en dé-commentant les lignes 40 à 70.

Attention, si ces lignes sont dé-commenter lors du dry-run, le job sera en erreur sur le cluster, étant donnée que le snakemake ne ce sera pas exécuter, il n'y aura aucune matrice à sommer.

TL;DR : Comment lancer une analyse sur Genologin :

- NfCore-hic :

1. Télécharger les reads au format fastq.gz dans des répertoires distincts `data/reads/Protocol`
2. Renseigner les champs du fichier `nfcore-hic/nfcore-hic.sh` .
3. Exécuter `sh nfcore-hic/nfcore-hic.sh` .
4. Soumettre le job `sbatch nfcore-hic/Protocol/Protocol.sh`

Puis une fois le nextflow terminé ;

- hicExplorer :

4. Renseigner les champs du fichier `hicexplorer/config.yaml`
5. Renseigner le chemin des matrices dans `hicexplorer/make_sampled_protocol.py`
6. Réaliser un dry-run en exécutant : `sankemake -j 1 -p -n` dans le dossier `hicexplorer`
7. Enlever le `-n` à la ligne 33 de `hicexplorer.sh`
8. Soumettre l'analyse `sbatch hicexplorer.sh`