

1. Сравнение реплик инструментом hicrep.

```
hicrep /home/kononkova/hic_data_hse/lib1.dm3.mapq_30.1000.mcool
/home/kononkova/hic_data_hse/lib2.dm3.mapq_30.1000.mcool outputSCC_lib.txt --binSize 10000
--h 1 --dBPMMax 5000000

hicrep /home/kononkova/hic_data_hse/Kc167_rep1.dm3.mapq_30.1000.mcool
/home/kononkova/hic_data_hse/Kc167_rep2.dm3.mapq_30.1000.mcool outputSCC_Kc167.txt --
binSize 10000 --h 1 --dBPMMax 5000000

hicrep /home/kononkova/hic_data_hse/lib1.dm3.mapq_30.1000.mcool
/home/kononkova/hic_data_hse/Kc167_rep2.dm3.mapq_30.1000.mcool outputSCC_lib_Kc_1.txt --
binSize 10000 --h 1 --dBPMMax 5000000

hicrep /home/kononkova/hic_data_hse/Kc167_rep1.dm3.mapq_30.1000.mcool
/home/kononkova/hic_data_hse/lib2.dm3.mapq_30.1000.mcool outputSCC_lib_Kc_2.txt --binSize
10000 --h 1 --dBPMMax 5000000

hicrep /home/kononkova/hic_data_hse/lib1.dm3.mapq_30.1000.mcool
/home/kononkova/hic_data_hse/Kc167_rep1.dm3.mapq_30.1000.mcool outputSCC_lib_Kc_3.txt --
binSize 10000 --h 1 --dBPMMax 5000000

hicrep /home/kononkova/hic_data_hse/lib2.dm3.mapq_30.1000.mcool
/home/kononkova/hic_data_hse/Kc167_rep2.dm3.mapq_30.1000.mcool outputSCC_lib_Kc_4.txt --
binSize 10000 --h 1 --dBPMMax 5000000
```

Для расчётов я воспользуюсь python.

```
B [1]: import pandas as pd
import numpy as np
from scipy.cluster import hierarchy
import matplotlib.pyplot as plt
```

```
B [2]: scc_list = ('outputSCC_lib.txt',
                  'outputSCC_Kc167.txt',
                  'outputSCC_lib_Kc_1.txt',
                  'outputSCC_lib_Kc_2.txt',
                  'outputSCC_lib_Kc_3.txt',
                  'outputSCC_lib_Kc_4.txt')

scc_dist = []

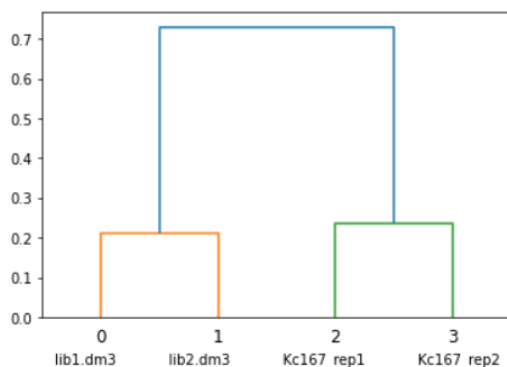
for i in range(len(scc_list)):
    with open(scc_list[i], 'r') as f:
        lib = f.readlines()
        scc = [float(x.strip()) for x in lib[2:-1]]
        scc = [i for x, i in enumerate(scc) if x in [x - 1 for x in [1, 3, 5, 7, 9, 12]]]
        scc_dist.append(np.mean(scc))

print(scc_dist)

[0.85045980851829, 0.8368286767714815, 0.5409096974046415, 0.5684901964584212, 0.5727378813725648, 0.5404668147418871]
```

```
B [3]: names = ('lib1.dm3', 'lib2.dm3', 'Kc167_rep1', 'Kc167_rep2')
dist_matrix = np.array([[1, scc_dist[0], scc_dist[4], scc_dist[2]],
                        [scc_dist[0], 1, scc_dist[3], scc_dist[5]],
                        [scc_dist[4], scc_dist[3], 1, scc_dist[1]],
                        [scc_dist[2], scc_dist[5], scc_dist[1], 1]])

dm = hierarchy.linkage(dist_matrix, 'single')
dn = hierarchy.dendrogram(dm)
plt.xlabel("".join(names)) # пришлось сделать так, подписи не хотели рисоваться красиво
plt.show()
```



Чего и следовало ожидать, реплики нервной и эмбриональной ткани разделились попарно.

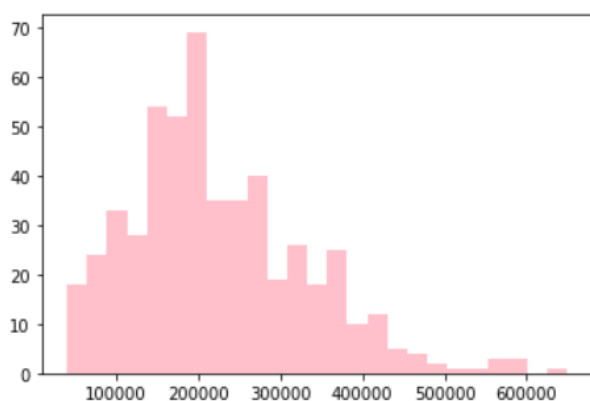
2. TADs calling.

Запускаю hicFindTADs со стандартными параметрами.

```
hicFindTADs -m /home/kononkova/hic_data_hse/lib_1_and_2.dm3.mapq_30.1000.mcool::/  
resolutions/10000 --outPrefix TADs --correctForMultipleTesting fdr --chromosomes chr2L  
chr2R chr3L chr3R chrX
```

```
: tads_default = pd.read_csv('tads/TADs_domains.bed', sep='\t', header=None)  
tads_default_len = [x - y for x, y in zip(tads_default[2].tolist(), tads_default[1].tolist())]  
print('mean length of TADs', np.mean(tads_default_len, dtype=int))  
plt.hist(tads_default_len, bins=25, color='pink')  
plt.show()
```

mean length of TADs 224575

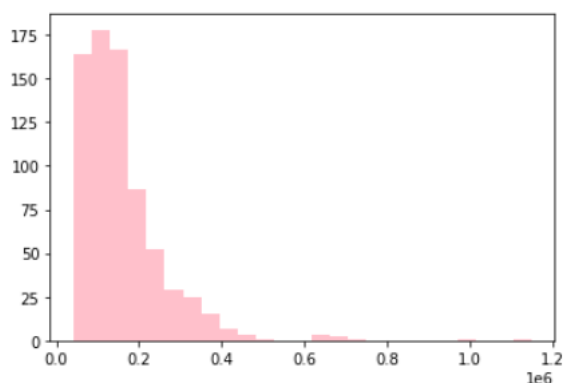


Я сделала несколько запусков с разными значениями `--minDepth` (30k, 100k, 200k) и `--maxDepth` (50k, 100k, 250k, 500k, 1M). Команды аналогичны этой:

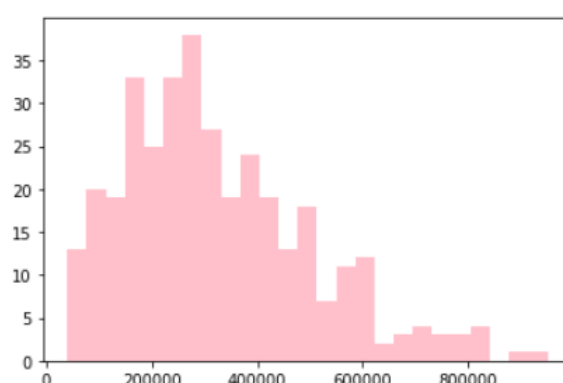
```
hicFindTADs -m /home/kononkova/hic_data_hse/lib_1_and_2.dm3.mapq_30.1000.mcool::/  
resolutions/10000 --outPrefix TAD_30k_100k --correctForMultipleTesting fdr --chromosomes  
chr2L chr2R chr3L chr3R chrX --minDepth 30000 --maxDepth 100000
```

При небольших значениях `minDepth` средняя длина ТАДов уменьшается и маленьких становится сильно больше. Но разброс значений при этом большой. С увеличением `minDepth` ТАДы становятся больше и разброс уменьшается.

`--minDepth 30k --maxDepth 100k`
mean length of TADs 157921



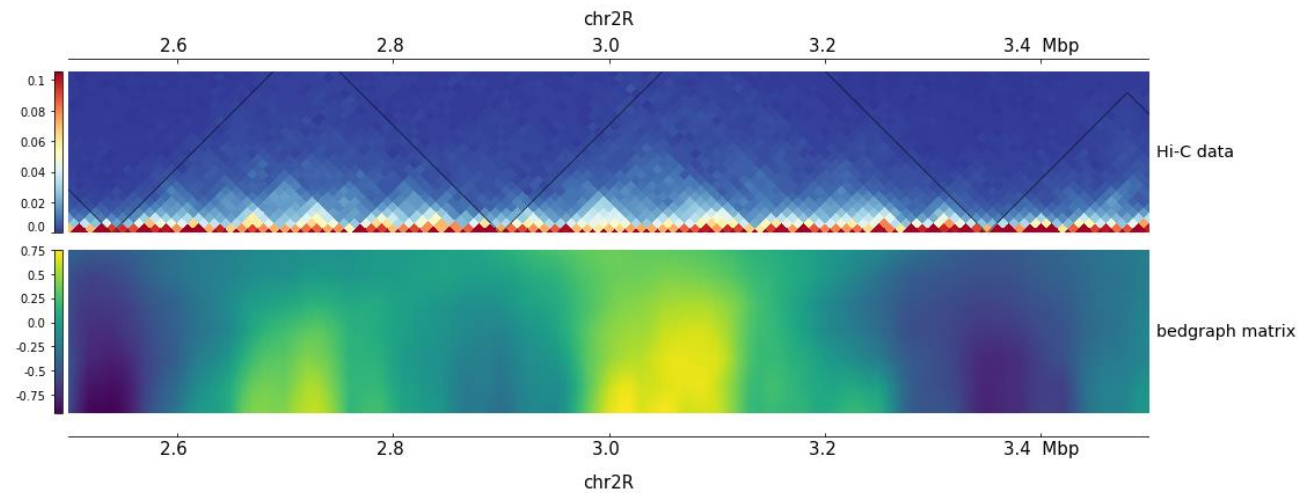
`--minDepth 200k --maxDepth 1M`
mean length of TADs 325071



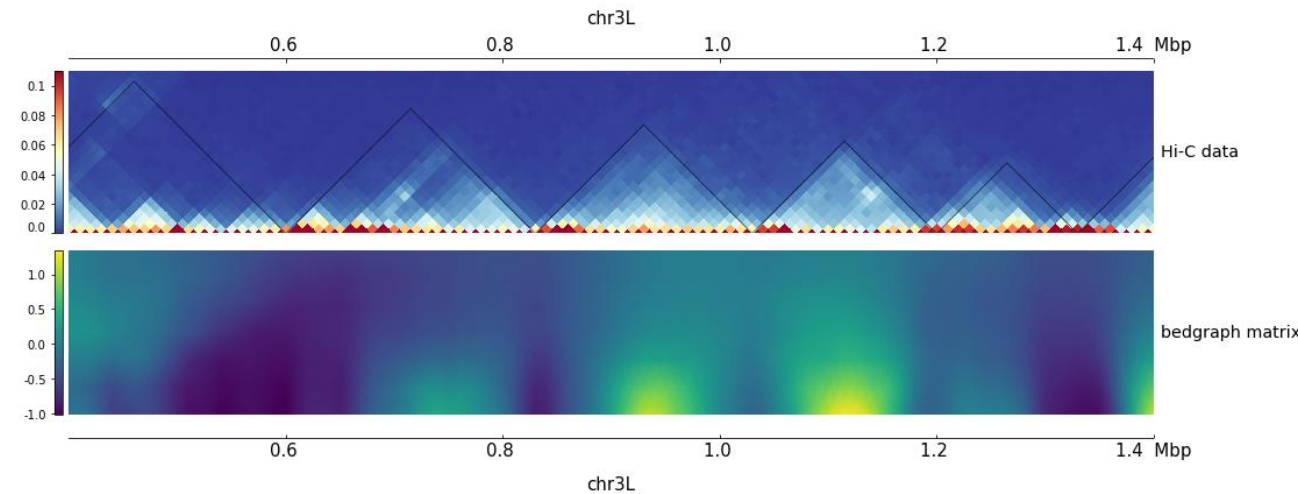
3. Визуализация.

Я выбрала два региона на разных хромосомах. И попробовала визуализировать ещё bedgraph matrix, получившиеся после hicFindTADs.

```
hicPlotTADs --tracks tracks.ini --region chr2R:2500100-3500100 -o drosophila_tads_2R_bm.png
```



```
hicPlotTADs --tracks tracks.ini --region chr3L:400050-1400050 -o drosophila_tads_3L_bm.png
```



4. Взаимосвязь между границами ТАДов и экспрессией.

Информацию о генах *Drosophila melanogaster* я скачала с Ensembl. Далее также считала в python.

BLAST/BLAT | VEP | Tools | BioMart | Downloads | Help & Docs | Blog

New | Count | Results

URL | XML | Perl | Help

Dataset

Drosophila melanogaster genes (BDGP6.28)

Filters

[None selected]

Attributes

Gene stable ID

Gene start (bp)

Gene end (bp)

Strand

Dataset

[None Selected]

Export all results to

File

CSV

☒ Unique results only

Go

Email notification to

View

All

rows as

HTML

☒ Unique results only

Gene stable ID	Gene start (bp)	Gene end (bp)	Strand
FBgn0031208	7529	9484	1
FBgn0002121	9839	21376	-1
FBgn0031209	21823	25155	-1
FBgn0263584	21952	24237	1
FBgn0051973	25402	65404	-1
FBgn0267987	54817	55767	1
FBgn0266878	65999	66242	1
FBgn0266879	66318	66524	1
FBgn0067779	66482	71390	1
FBgn0266322	71039	73836	-1

Координаты в файле *boundaries.bed* содержат координаты границ ТАДов. В то время как *domains.bed* содержит координаты самих ТАДов. Я решила, что если начало или конец гена попадает в промежуток между началом и концом границы ТАДа, будем считать, что он относится к границе. Если не попадает в этот промежуток – ген находится в ТАДе. Для поиска по всем генам получилось, что больше генов находится в граничных регионах.

```
gene_start = drosophila['Gene start (bp)'].tolist()
gene_end = drosophila['Gene end (bp)'].tolist()

i = 0 # итератор генов
b = 0 # попадает на границу
x = 0 # падает внутрь ТАДа
for n in range(len(TADs_boundaries)):
    while gene_end[i] < TADs_boundaries[1][n]:
        i += 1
    if TADs_boundaries[1][n] < gene_start[i] < TADs_boundaries[2][n] or \
    TADs_boundaries[1][n] < gene_end[i] < TADs_boundaries[2][n]:
        while TADs_boundaries[1][n] < gene_start[i] < TADs_boundaries[2][n] or \
        TADs_boundaries[1][n] < gene_end[i] < TADs_boundaries[2][n]:
            b += 1
            i += 1
    elif TADs_boundaries[2][n] < gene_start[i] < TADs_boundaries[1][n+1] or \
    TADs_boundaries[2][n] < gene_end[i] < TADs_boundaries[1][n+1]:
        if n < len(TADs_boundaries):
            while TADs_boundaries[2][n] < gene_start[i] < TADs_boundaries[1][n+1] or \
            TADs_boundaries[2][n] < gene_end[i] < TADs_boundaries[1][n+1]:
                x += 1
                i += 1
    elif TADs_boundaries[2][n] < gene_start[i]:
        while n < len(TADs_boundaries):
            if TADs_boundaries[2][n] > gene_start[i]:
                break
            n += 1
        continue
    if n == len(TADs_boundaries):
        break

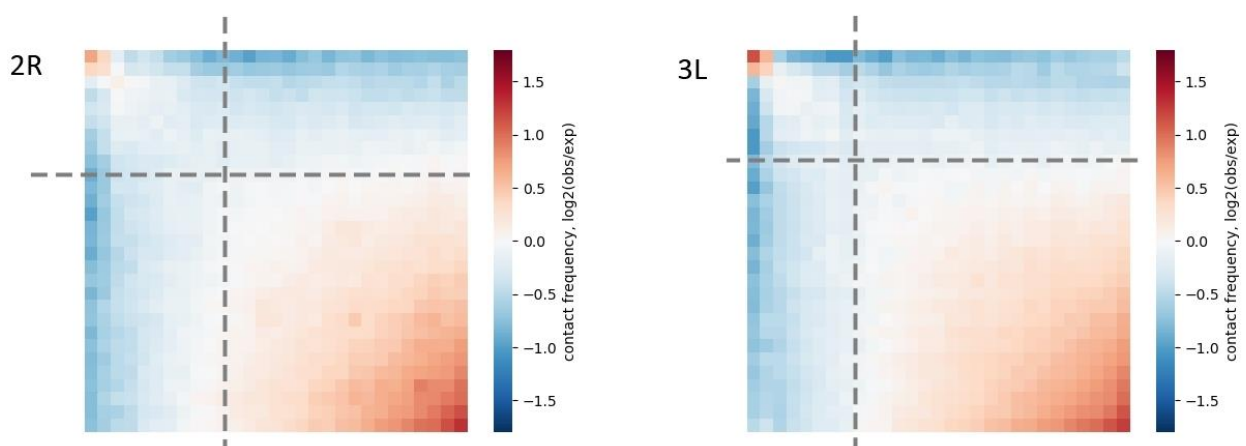
print('in boundaries:', b)
print('in TADs:', x)

in boundaries: 447
in TADs: 246
```

Аналогичный алгоритм я повторила на генах с максимальной экспрессией. В скачанном мной датасете с Ensembl нашлось 340 и 400 генов. И получилось, что на границы попадает 19 генов, а в ТАДы – 56. Следовательно, в хорошо экспрессированные гены чаще находятся в ТАДах (jupyter notebook прилагаю).

5. Компартменты.

Saddle plot для хромосом 2R и 3L. Эти графики показывают частоту контактов по длине хромосомы. У обеих хромосом картина получается очень схожая: одно плечо хромосомы сильно меньше другого.



6. Получение Hi-C карт из fastq-файлов.

Выход моей программы.

```
(base) anastasia_rastvorova@hpc09:~/hi-c/dstlr$ nextflow distiller.nf -params-file ./project.yml
N E X T F L O W ~ version 20.10.0
Launching `distiller.nf` [amazing_mercator] - revision: af64431922
executor > local (14)
[-] process > download_truncate_chunk_fastqs -
[skipped] process > local_truncate_chunk_fastqs (library:replica_1 run:lane1) [100%] 2 of 2, stored: 2 ✓
[-] process > fastqc -
[fa/45bb23] process > map_parse_sort_chunks (library:replica_2 run:lane1 chunk:0) [100%] 2 of 2 ✓
[59/b77b68] process > merge_dedup_splitbam (library:replica_1) [100%] 2 of 2 ✓
[51/de220c] process > bin_zoom_library_pairs (library:replica_1 filter:no_filter) [100%] 4 of 4 ✓
[b1/2ca64b] process > merge_zoom_library_group_coolers (library_group:all filter:mapq_30) [100%] 4 of 4 ✓
[09/20a595] process > merge_stats_libraries_into_groups (library_group:all) [100%] 2 of 2 ✓
Completed at: 18-Dec-2020 22:23:45
Duration : 13m 46s
CPU hours : 1.8
Succeeded : 14
```

Посмотрим, что получилось.

```
cooler info results/cooler_library_group/BG3.dm3.no_filter.1000.mcool::/resolutions/10000
```

```
cooler info results/cooler_library_group/BG3.dm3.mapq_30.1000.mcool::/resolutions/10000
```

```
(base) anastasia_rastvorova@hpc09:~/hi-c/dstlr$ cooler info results/cooler_library_group/BG3.dm3.no_filter.1000.mcool::/resolutions/10000
{
  "bin-size": 10000,
  "bin-type": "fixed",
  "creation-date": "2020-12-18T22:19:53.145351",
  "format": "HDF5::Cooler",
  "format-url": "https://github.com/mirnylab/cooler",
  "format-version": 3,
  "generated-by": "cooler-0.8.10",
  "genome-assembly": "unknown",
  "nbins": 16880,
  "nchroms": 15,
  "nnz": 382553,
  "storage-mode": "symmetric-upper",
  "sum": 1055863
}
(base) anastasia_rastvorova@hpc09:~/hi-c/dstlr$ cooler info results/cooler_library_group/BG3.dm3.mapq_30.1000.mcool::/resolutions/10000
{
  "bin-size": 10000,
  "bin-type": "fixed",
  "creation-date": "2020-12-18T22:18:51.837230",
  "format": "HDF5::Cooler",
  "format-url": "https://github.com/mirnylab/cooler",
  "format-version": 3,
  "generated-by": "cooler-0.8.10",
  "genome-assembly": "unknown",
  "nbins": 16880,
  "nchroms": 15,
  "nnz": 366344,
  "storage-mode": "symmetric-upper",
  "sum": 1022730
}
```

Без фильтра (1055863), как понимаю, контактов получается больше, чем с фильтром (1022730).

Здесь же прилагаю, что у меня написано в *project.yml* (файл прилагаю).

```

1 input:
2   raw_reads_paths:
3     replica_1:
4       lane1:
5         - /home/kononkova/hic_data_hse/SRR8195120_1.fastq.gz
6         - /home/kononkova/hic_data_hse/SRR8195120_2.fastq.gz
7     replica_2:
8       lane1:
9         - /home/kononkova/hic_data_hse/SRR8195121_1.fastq.gz
10        - /home/kononkova/hic_data_hse/SRR8195121_2.fastq.gz
11
12   library_groups:
13     BG3:
14       - replica_1
15       - replica_2
16
17   truncate_fastq_reads: 1000000
18
19   genome:
20     assembly_name: 'dm3'
21     bwa_index_wildcard_path: '/mnt/local/vse2020/home/anastasia_rastvorova/hi-c/dm3.fa.gz.*'
22     chrom_sizes_path: '/mnt/local/vse2020/home/anastasia_rastvorova/hi-c/dm3.chrom.sizes.txt'
23
24   do_fastqc: False
25
26   map:
27     chunksize: 0
28     mapping_options: ''
29     trim_options: ''
30     use_custom_split: true
31
32   parse:
33     make_pairsam: False
34     drop_seq: False
35     drop_readid: True
36     keep_unparsed_bams: False
37     parsing_options: '--add-columns mapq --walks-policy mask'
38
39   dedup:
40     max_mismatch_bp: 1
41
42   bin:
43     resolutions:
44       - 1000000
45       - 500000
46       - 250000
47       - 100000
48       - 50000
49       - 25000
50       - 10000
51       - 5000
52       - 2000
53       - 1000
54     balance: true
55     filters:
56       no_filter: ''
57       mapq_30: '(mapq1>=30) and (mapq2>=30)'
58
59   output:
60     dirs:
61       processed_fastqs: 'results/processed_fastqs/'
62       mapped_parsed_sorted_chunks: 'results/mapped_parsed_sorted_chunks/'
63       fastqc: 'results/fastqc/'
64       pairs_library: 'results/pairs_library/'
65       coolers_library: 'results/coolers_library/'
66       coolers_library_group: 'results/coolers_library_group/'
67       stats_library_group: 'results/stats_library_group/'
68

```