Homework 1

Name: Aris Podotas

University: National and Kapodistrian University of Athens

Program: Data Science and Informaion Technologies

Specialization: Bioinformatics - Biomedical Data

Lesson: Introduction To Bioinformatics

User: user03

Date: November 2024

Contents

1	Ana	alysis	
	1.1	find a read with unmapped pair	
	1.2	find a read that has 2 mismatches	
	1.3	find a read-pair with pair orientation: F2R1	
	1.4	find a gene that agrees in 6 splice-sites with its annotation	
	1.5	find a location where less than 96at least 2 other nucleotides occur	
	1.6	find a gene where the peak with the highest number of reads has at least	
		3 times the reads of the highest peak in a different sample (with non-zero reads for that gene)	
2			
3			
	3.1	Right reads	
	3.2	Left reads	
4			
	4.1	Attempt 1:	
	4.2	Attempt 2:	
	4.3	Attempt 3:	
5	Opt	Optional task	



Note: Files will be within the user03 folder in the virtual machine. **Note:** Commands will not be explicitly written here.

1 Analysis

Note: a intron has two splice sites, the 5' and 3' splice site.

1.1 find a read with unmapped pair

genome:173,708-173,775

1.2 find a read that has 2 mismatches

genome:136,210-136,277

1.3 find a read-pair with pair orientation: F2R1

genome:313,112-313,179

- 1.4 find a gene that agrees in 6 splice-sites with its annotation genome:329662-331216
- 1.5 find a location where less than 96at least 2 other nucleotides occur.

genome:270,404

1.6 find a gene where the peak with the highest number of reads has at least 3 times the reads of the highest peak in a different sample (with non-zero reads for that gene)

genome:2670-3431

2

After running the command (the command itself is contained in the file) and looking in the "cuffcmp.stats" file we have:

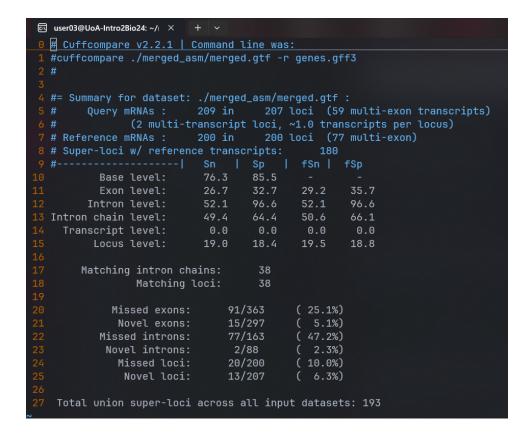


Figure 1: Cuffcompare: cuffcmp.stats file contents

3

Showing the contents of all the files under the align_summary.txt label:

```
Left reads:
          Input
                         101577
                          92230 (90.8% of input)
           Mapped
                           14 ( 0.0%) have multiple alignments (0 have >20)
Right reads:
          Input
                          95250 (93.8% of input)
           Mapped
                             13 ( 0.0%) have multiple alignments (0 have >20)
92.3% overall read mapping rate.
Aligned pairs:
                   88881
                      12 ( 0.0%) have multiple
                      31 ( 0.0%) are discordant
                                                     ments
87.5% concordant pair
                          nment rate.
```

Figure 2: Tophat: df summary statistics

```
Left reads:
          Input
                         107676
          Mapped
                            19 ( 0.0%) have multiple alignments (0 have >20)
Right reads:
          Input
                        107676
          Mapped :
                        100451 (93.3% of input)
                         18 ( 0.0%) have multiple alignments (0 have >20)
92.7% overall read mapping rate.
                  95068
Aligned pairs:
                     18 ( 0.0%) have multiple
    of these:
                                                   ments
                        ( 0.0%) are discordant
                                                    ments
88.3% concordant pair
                          ment rate.
```

Figure 3: Tophat: hs summary statistics

Figure 4: Tophat: log summary statistics

```
Left reads:
                          99328
           Mapped:
                          91593 (92.2% of input)
                            19 ( 0.0%) have multiple alignments (0 have >20)
            of these:
Right reads:
          Input
          Mapped :
                          94339 (95.0% of input)
                           19 ( 0.0%) have multiple alignments (0 have >20)
           of these:
93.6% overall read mapping rate.
Aligned pairs:
                      17 ( 0.0%) have multiple
    of these:
                                                    ments
                      14 ( 0.0%) are discordant
89.6% concordant pair
                           ment rate.
```

Figure 5: Tophat: plat summary statistics

3.1 Right reads

The final right reads average is:

$$\frac{93.8 + 93.3 + 93.3 + 95.0}{4} = 93.85\%$$



3.2 Left reads

The final left reads average is:

$$\frac{90.8 + 92.1 + 92.2 + 92.2}{4} = 91.825\%$$

4

Note: The files mentioned will be within user03 in the virtual machine

Note: Input functionality at each step will be

Note: Because all answers will be yes for this exercise as was said in the requirements.

Note: The files that call each other need *chmod* to be run on them before we can use them since they dont have permission to run other files yet.

Multiple attempts were made and documented (kept within the vm). Each one had a layer of progression, for example the first attempt had one file that ran on one thread not even qualifying for the solution. These are the files:

4.1 Attempt 1:

File: script.sh

```
#!/bin/bash
2
     main () {
3
            bowtie-build genome.fa genome
4
                                          20
            tophat -I 1000 -i
                                                 --bowtie1
    library-type fr-firststrand -o
                                  tophat.Sp_ds.dir
                                                              genome
     Sp_ds.left.fq Sp_ds.right.fq
                   tophat.Sp_ds.dir/accepted_hits.bam
            mν
    Sp_ds.dir/Sp_ds.bam
                           index tophat.Sp_ds.dir/Sp_ds.bam
            samtools
8
            cufflinks
                           --overlap-radius 1 --library-type fr-
                                            tophat.Sp_ds.dir/Sp_ds.
    firststrand -o cufflinks.Sp_ds.dir
    bam
     mv cufflinks.Sp_ds.dir/transcripts.gtf cufflinks.Sp_ds.dir/Sp_ds.
    transcripts.gtf
            tophat -I
                           1000
                                  -i
                                          20
                                                 --bowtie1
    library-type fr-firststrand -o
                                       tophat.Sp_hs.dir
                                                              genome
      Sp_hs.left.fq Sp_hs.right.fq
                    tophat.Sp_hs.dir/accepted_hits.bam
    Sp_hs.dir/Sp_hs.bam
                           index tophat.Sp_hs.dir/Sp_hs.bam
            samtools
                         --overlap-radius 1
            cufflinks
                                                         --library-
    type fr-firststrand -o cufflinks.Sp_hs.dir tophat.Sp_hs.
    dir/Sp_hs.bam
                   cufflinks.Sp_hs.dir/transcripts.gtf
                                                         cufflinks.
    Sp_hs.dir/Sp_hs.transcripts.gtf
                                          20
                                                 --bowtie1
            tophat -I 1000 -i
    library-type fr-firststrand -o
                                       tophat.Sp_log.dir
                                                              genome
      Sp_log.left.fq Sp_log.right.fq
```



```
tophat.Sp_log.dir/accepted_hits.bam
                                                                tophat.
     Sp_log.dir/Sp_log.bam
                               index
                                       tophat.Sp_log.dir/Sp_log.bam
              samtools
17
                               --overlap-radius 1
              cufflinks
                                                                --library-
18
     type fr-firststrand -o
                                   cufflinks.Sp_log.dir
                                                            tophat.Sp_log.
     dir/Sp_log.bam
                      cufflinks.Sp_log.dir/transcripts.gtf
                                                                cufflinks.
19
     Sp_log.dir/Sp_log.transcripts.gtf
              tophat -I
                              1000
                                               20
20
     library-type fr-firststrand -o
                                            tophat.Sp_plat.dir
                      tophat.Sp_plat.dir/accepted_hits.bam
                                                                tophat.
21
              mν
     Sp_plat.dir/Sp_plat.bam
                               index
                                       tophat.Sp_plat.dir/Sp_plat.bam
22
              samtools
              cufflinks --overlap-radius 1
                                             --library-type fr-
23
                          cufflinks.Sp_plat.dir tophat.Sp_plat.dir/
     firststrand -o
     Sp_plat.bam
                      cufflinks.Sp_plat.dir/transcripts.gtf
24
              mν
     Sp_plat.dir/Sp_plat.transcripts.gtf
                      cufflinks.Sp_ds.dir/Sp_ds.transcripts.gtf
                                                                        >>
          assemblies.txt
                      cufflinks.Sp_hs.dir/Sp_hs.transcripts.gtf
26
          assemblies.txt
                      cufflinks.Sp_log.dir/Sp_log.transcripts.gtf
                                                                        >>
              echo
          assemblies.txt
              echo
                      cufflinks.Sp_plat.dir/Sp_plat.transcripts.gtf
28
          assemblies.txt
                      assemblies.txt
              cat
              cuffmerge
                                       genome.fa
                                                       assemblies.txt
30
                     -Xmx2G -jar
                                      /Users/bhaas/IGV/current//igv.jar
31
                  'pwd'/genome.fa 'pwd'/merged_asm/merged.gtf, 'pwd'/genes
     .bed, 'pwd'/tophat.Sp_ds.dir/Sp_ds.bam, 'pwd'/tophat.Sp_hs.dir/Sp_hs.
     bam, 'pwd'/tophat.Sp_log.dir/Sp_log.bam, 'pwd'/tophat.Sp_plat.dir/
     Sp_plat.bam
              cuffdiff
                               --library-type fr-firststrand
39
                                                                        -0
          diff_out
                          -b
                                   genome.fa
                                                   - L
                                                            Sp_ds,Sp_hs,
                               merged_asm/merged.gtf
                                                         tophat.Sp_ds.dir/
     Sp_log,Sp_plat
                         -u
                    tophat.Sp_hs.dir/Sp_hs.bam
                                                     tophat.Sp_log.dir/
     Sp_ds.bam
                   tophat.Sp_plat.dir/Sp_plat.bam
     Sp_log.bam
                      diff_out/gene_exp.diff
              head
34
              return 0
      }
35
36
37 main
```

This is a single file non working version that is literally the commands from the pdf file transferred to a bash file. This does not constitute a solution to the problem.

4.2 Attempt 2:

File: mthscript.sh

```
#!/bin/bash
main () {
```

Data Science and Information Technologies Master's National and Kapodistrian University of Athens

```
bowtie-build genome.fa genome -p 4
    tophat -p 4 -I 1000 -i 20 --bowtie1
library-type fr-firststrand -o tophat.Sp_ds.dir
                                                                 genome
      Sp_ds.left.fq Sp_ds.right.fq
            mv tophat.Sp_ds.dir/accepted_hits.bam tophat.
     Sp_ds.dir/Sp_ds.bam
    samtools -@ 4 index tophat.Sp_ds.dir/Sp_ds.bam cufflinks -p 4 --overlap-radius 1 --library-type fr-firststrand -o cufflinks.Sp_ds.dir tophat.Sp_ds.dir/Sp_ds.
            mv cufflinks.Sp_ds.dir/transcripts.gtf cufflinks.Sp_ds.dir
9
     /Sp_ds.transcripts.gtf
             tophat -p 4 -I 1000 -i
                                         20 --bowtie1
     library-type fr-firststrand -o tophat.Sp_hs.dir
                                                                 genome
       Sp_hs.left.fq Sp_hs.right.fq
                     tophat.Sp_hs.dir/accepted_hits.bam tophat.
     Sp_hs.dir/Sp_hs.bam
     samtools -0 4 index tophat.Sp_hs.dir/Sp_hs.bam
cufflinks -p 4 --overlap-radius 1 --library-
type fr-firststrand -o cufflinks.Sp_hs.dir tophat.Sp_hs.
                            -@ 4 index tophat.Sp_hs.dir/Sp_hs.bam
12
     dir/Sp_hs.bam
                    cufflinks.Sp_hs.dir/transcripts.gtf
                                                           cufflinks.
14
     Sp_hs.dir/Sp_hs.transcripts.gtf
                                         20 --bowtie1
             tophat -p 4 -I 1000 -i
     library-type fr-firststrand -o tophat.Sp_log.dir
                                                                 genome
     Sp_log.left.fq Sp_log.right.fq
                    tophat.Sp_log.dir/accepted_hits.bam tophat.
     Sp_log.dir/Sp_log.bam
                           -@ 4 index tophat.Sp_log.dir/Sp_log.
            samtools
     bam
             cufflinks -p 4 --overlap-radius 1
     type fr-firststrand -o cufflinks.Sp_log.dir tophat.Sp_log.
     dir/Sp_log.bam
             mv
                    cufflinks.Sp_log.dir/transcripts.gtf cufflinks.
     Sp_log.dir/Sp_log.transcripts.gtf
                                        20 --bowtie1
     20
             mv tophat.Sp_plat.dir/accepted_hits.bam
                                                           tophat.
     Sp_plat.dir/Sp_plat.bam
                             -@ 4 index tophat.Sp_plat.dir/Sp_plat.
             samtools
     bam
             cufflinks -p 4 --overlap-radius 1
                                                    --library-<mark>type</mark> fr-
     firststrand -o cufflinks.Sp_plat.dir tophat.Sp_plat.dir/
     Sp_plat.bam
             mv cufflinks.Sp_plat.dir/transcripts.gtf cufflinks.
     {\tt Sp\_plat.dir/Sp\_plat.transcripts.gtf}
             echo cufflinks.Sp_ds.dir/Sp_ds.transcripts.gtf
         assemblies.txt
             echo cufflinks.Sp_hs.dir/Sp_hs.transcripts.gtf
                                                                    >>
         assemblies.txt
             echo cufflinks.Sp_log.dir/Sp_log.transcripts.gtf
          assemblies.txt
             echo cufflinks.Sp_plat.dir/Sp_plat.transcripts.gtf
          assemblies.txt
             cat assemblies.txt
          cuffmerge -p 4 -s genome.fa assemblies.txt
30
```



```
java -Xmx2G -jar /Users/bhaas/IGV/current//igv.jar
                'pwd'/genome.fa 'pwd'/merged_asm/merged.gtf, 'pwd'/genes
     .bed, 'pwd'/tophat.Sp_ds.dir/Sp_ds.bam, 'pwd'/tophat.Sp_hs.dir/Sp_hs.
     bam, 'pwd'/tophat.Sp_log.dir/Sp_log.bam, 'pwd'/tophat.Sp_plat.dir/
     Sp_plat.bam
                            --library-type fr-firststrand
             cuffdiff
32
                        -b
                                                      Sp_ds, Sp_hs,
         diff_out
                              genome.fa -L
     Sp_log,Sp_plat -u merged_asm/merged.gtf tophat.Sp_ds.dir/
     Sp_ds.bam tophat.Sp_hs.dir/Sp_hs.bam tophat.Sp_log.dir/
     Sp_log.bam tophat.Sp_plat.dir/Sp_plat.bam
                    diff_out/gene_exp.diff
             head
33
             return 0
35
     }
36
37
     main
```

This is a modified version of the script that runs on 4 threads, there is still some amount of repetition in the code that should be abstracted away. It is a single file.

4.3 Attempt 3:

Files: executer.sh, abstract.sh (in that order)

```
#!/bin/bash
2
      main() {
              runs=('ds' 'hs' 'log' 'plat')
              bowtie-build genome.fa genome -p 4
              for data in "${runs[@]}"
                      echo "Running
     $data"
                      ./abstract.sh $data
10
                      -Xmx2G -jar /Users/bhaas/IGV/current//igv.jar
              'pwd'/genome.fa 'pwd'/merged_asm/merged.gtf, 'pwd'/genes
     .bed, 'pwd'/tophat.Sp_ds.dir/Sp_ds.bam, 'pwd'/tophat.Sp_hs.dir/Sp_hs.
     bam, 'pwd'/tophat.Sp_log.dir/Sp_log.bam, 'pwd'/tophat.Sp_plat.dir/
     Sp_plat.bam
              cuffdiff
                               --library-type fr-firststrand
                       -b genome.fa -L Sp_ds,Sp_hs,
-u merged_asm/merged.gtf tophat.Sp_ds.dir/
          diff_out
     Sp_log,Sp_plat
     Sp_ds.bam
                   tophat.Sp_hs.dir/Sp_hs.bam
                                                    tophat.Sp_log.dir/
     Sp_log.bam tophat.Sp_plat.dir/Sp_plat.bam
              head diff_out/gene_exp.diff
      }
14
      main
```

1 #!/bin/bash
2



```
tophat -p 4 -I 1000 -i 20 --bowtie1 --library-type
    fr-firststrand -o tophat.Sp_$1.dir genome Sp_ds.left.fq
    Sp_ds.right.fq

mv tophat.Sp_$1.dir/accepted_hits.bam tophat.Sp_$1.dir/Sp_$1.
    bam

samtools index tophat.Sp_$1.dir/Sp_$1.bam -@ 4

cufflinks -p 4 --overlap-radius 1 --library-type fr-firststrand
    -o cufflinks.Sp_$1.dir tophat.Sp_$1.dir/Sp_$1.bam

mv cufflinks.Sp_$1.dir/transcripts.gtf cufflinks.Sp_$1.dir/
    Sp_$1.transcripts.gtf

echo cufflinks.Sp_$1.dir/Sp_$1.transcripts.gtf >> assemblies.txt

cat assemblies.txt

cuffmerge -p 4 -s genome.fa assemblies.txt
```

These are two files that interplay running the commands with 4 threads and using a more efficient method for calling similar commands with slight name changes such that the lines of code actually written are minimized.

5 Optional task

Files: hsat1.sh, hsat2.sh

```
#!/bin/bash
2
3 main() {
         runs=('ds' 'hs' 'log' 'plat')
         hisat2-build genome.fa genome -p 4
         echo > assemblies.txt
6
         for data in "${runs[0]}"
                 echo "Running
     $data"
                 ./hsat2.sh $data
11
                       --library-type fr-firststrand
         cuffdiff
12
     hsat -b genome.fa -L Sp_ds,Sp_hs,Sp_log,Sp_plat -u
        merged_asm/merged.gtf Sp_ds.sorted.bam Sp_hs.sorted.bam
            Sp_log.sorted.bam
                                Sp_plat.sorted.bam
13 }
14
15 main
#!/bin/bash
3 hisat2 -p 4 --max-intronlen 1000 --fr --min-intronlen 20
     -x genome -1 Sp_$1.left.fq
                                          -2 Sp_$1.right.fq -S Sp_$1.
     sam
4 samtools view -bS Sp_$1.sam > Sp_$1.bam
5 samtools sort -o Sp_$1.sorted.bam Sp_$1.bam
6 samtools index Sp_$1.sorted.bam
7 cufflinks -p 4 --overlap-radius 1 --library-type fr-firststrand
-o hisat.cufflinks.Sp_$1.dir Sp_$1.sorted.bam
8 mv hisat.cufflinks.Sp_$1.dir/transcripts.gtf hisat.cufflinks
 .Sp_$1.dir/Sp_$1.transcripts.gtf
```



Data Science and Information Technologies Master's National and Kapodistrian University of Athens

```
echo hisat.cufflinks.Sp_$1.dir/Sp_$1.transcripts.gtf >>
    assemblies.txt

cat assemblies.txt
cuffmerge -p 4 -s genome.fa assemblies.txt
```

These resemble the executer.sh and abstract.sh because they are modified copies. the alterations to the flags and the commands came from reading the documentation of each of the tools and finding common flags with -h.

For the difference between the hisat2 and tophat.

Which essentially means that hisat 2 produces more statistically significant results, we can expect that this is due to the different alignment method it uses. At face value it means that hisat2 identifies up regulated and down regulated genes but tophat does not.

Thus concludes this ITBI exercise