

Log: Genome Annotation Chrysomya putoria

tags: `Genome Annotation` `Cput`

Table of Contents

- [Log: Genome Annotation Chrysomya putoria](#)
 - [Table of Contents](#)
 - [List of software \(and where they are located in Rosalind\):](#)
 - [Genome annotation workflow:](#)
- [Genome sequencing](#)
- [0- Copy of the genome](#)
- [1-BUSCO first](#)
- [2- QUAST](#)
- [3- Mitochondrial genome](#)
 - [Blastn](#)
 - [MITOS Web Server](#)
 - [TEST on pacbio genome before assembly](#)
- [4- RepeatModeler](#)
 - [Installing with docker](#)
 - [Creating database](#)
 - [Running](#)
 - [Copying the results](#)
 - [Test - LTR Harvest](#)
- [5- RepeatMasker](#)
 - [2.5.1 BUSCO](#)
 - [2.5.2 QUAST](#)
- [6- RNA-seq](#)
 - [Quality control of raw reads](#)
 - [Trimming](#)
 - [Quality control of trimmed reads](#)

- Observation
- 6.1 - New RNA sequencing
 - Quality control of raw reads
 - Trimming
 - Quality control of trimmed reads
 - Assembling the transcriptome ("old + new" RNA-seq)
 - BUSCO of the new transcriptome
- 7- Trinity
 - Moving trimmed reads to Darwin
 - Assembling the transcriptome
 - BUSCO (transcriptome quality)
- 8- STAR - RNAseq alignment
 - index
 - STAR alignment
- 9- Braker3
 - BRAKER3 second run
 - Busco evaluation with protein sequences from restart Braker run
- 11- EnTAP
- 12-Final files
- 14 - Submission to NCBI

List of software (and where they are located in Rosalind):

- QUAST (/usr/local/bin/quast.py)
- RepeatModuler (/RepeatModeler-2.0.4)
- RepeatMasker(/dados/home/pedro/Programs/RepeatMasker/RepeatMasker)
- STAR (/dados/home/bruno/anaconda3/bin/STAR)
- samtools (/dados/home/bruno/anaconda3/bin/samtools
/dados/home/bruno/anaconda3/bin/samtools.pl)
- Braker (docker image ID: 7772eca57cee)
 - <https://hub.docker.com/r/teambraker/braker3>
- TSEBRA (/)
 - <https://github.com/Gaius-Augustus/TSEBRA>

- AUGUSTUS (docker image ID: c0dfd27799fc)
 - <https://github.com/Gaius-Augustus/Augustus>
- busco (/usr/local/bin/busco)
- gFACS (/)
 - <https://gitlab.com/PlantGenomicsLab/gFACS>
- EnTAP (/)
 - <https://entap.readthedocs.io/en/v0.8.0-beta/introduction.html>

Installation for every program steps can be found in the *Chrysomya megacephala* log (https://hackmd.io/ziOztK1MQZecVm0_qDAr1g)

Genome annotation workflow:

- ☒ 1- BUSCO
- ☒ 2- QUAST
- ☒ 3- Mitochondrial genome
- ☒ 4- RepeatModeler
- ☐ 5- RepeatMasker
- ☒ 6- RNA-seq
- ☒ 7- Trinity
- ☒ 8- STAR
- ☐ 9- Braker

https://github.com/CBC-UCONN/Genome_Assembly
<https://github.com/CBC-UCONN/Structural-Annotation>

Genome sequencing

We sequenced a pool of males at Dovetail (Pac-Bio).

ANOTAR AQUI DADOS DO SEQUENCIAMENTO

Number of reads 2639925

Coverage (x) 58

HPA Length (bp) 1287183393

HPA N50 (bp) 1685985

HPA N90 (bp) 71006

HPA L50 168

HPA L90 1953

FA Length (bp) 587066816

FA N50 (bp) 4322519

FA N90 (bp) 590301

FA L50 45

FA L90 179

bp = Base pair; HPA = Hifiasm (Cheng et al., 2021) primary assembly; FA = Final assembly.
N50 = Sequence length of the smallest contig within those that sum up to 50% of the total genome's length; N90 = Sequence length of the smallest contig within those that sum up to 90% of the total genome's length; L50 = Smallest sequence number that together sum 50% of the total genome's length; L90 = Smallest sequence number that together sum 90% of the total genome's length.

0- Copy of the genome

We made an extra copy of the genome and the PacBio reads just to be safe.

```
1 # we are here
2 /home/blowflies/genome_annotation/cput
3
4 # new directory
5 mkdir 0-genome
6 cd 0-genome
7
8 cp /home/Reference_genomes/Cputoria/purged.fa .
9 cp /home/Reference_genomes/Cputoria/XDOVE_20221110_S64411e_PL100270437-1_D01.
```

1-BUSCO first

Busco version 5.3.2

```
1 mkdir /home/blowflies/genome_annotation/cput/1-busco_first
2 cd /home/blowflies/genome_annotation/cput/1-busco_first
3
4 #sudo docker pull ezlabgva/busco:v5.4.4_cv1
5
6 cp ../0-genome/cput.fa .
7
8 sudo docker run -u $(id -u) -v $(pwd):/busco_wd ezlabgva/busco:v5.4.4_cv1 bus
```

```
-----  
|Results from dataset diptera_odb10  
-----
```

```
|C:98.6%[S:93.3%,D:5.3%],F:0.3%,M:1.1%,n:3285  
|3241   Complete BUSCOs (C)  
|3066   Complete and single-copy BUSCOs (S)  
|175    Complete and duplicated BUSCOs (D)  
|10     Fragmented BUSCOs (F)  
|34     Missing BUSCOs (M)  
|3285   Total BUSCO groups searched  
-----
```

2- QUAST

```
1 | # we are here  
2 | /home/blowflies/genome_annotation/cput  
3 |  
4 | # new directory  
5 | mkdir 2-quast  
6 | cd 2-quast  
7 |  
8 | # running  
9 | sudo quast.py ../0-genome/purged.fa -t 10 --eukaryote --large --rna-finding
```

3- Mitochondrial genome

Cput mitochondrial genome on NCBI

<https://www.ncbi.nlm.nih.gov/nuccore/AF352790.1>

search terms: "chrysomya putoria[ORGN] AND mitochondrial genome"

AF352790.1

```
1 | esearch -db nuccore -query AF352790.1 | efetch -format fasta > AF352790.1.fa  
2 | cd home/blowflies/genome_annotation/calb/3-mitochondrial_genome  
3 | mkdir 1-blastn
```

We renamed the directory

```
1 | mv 3-mitochondrial_genome/ 3-mitochondrial_genome
```

Blastn

version 2.11

```

1 | cd home/blowflies/genome_annotation/cput/3-mitochondrial_genome/1-blastn
2 | #Making database using the genome
3 | makeblastdb -in ../AF352790.1.fa -dbtype nucl -out cput_mit_database
4 | #Running
5 | blastn -task blastn -evaluate 0.00001 -db ./cput_mit_database -query ../../0-genome

```

The best alignment was against scaffold 1221, which is a very big scaffold (0.5Mb). We isolated this scaffold and cut it in three pieces "before mitochondria" - "mitochondria" - "after mitochondria"

```

1 | cd /home/blowflies/genome_annotation/cput/0-genome
2 |
3 | # Remove scaffold from genome
4 | seqkit grep -v -p "ptg001221l" cput.fa > cput_N_genome.fa
5 |
6 | # Isolate scaffold
7 | seqkit grep -p "ptg001221l" cput.fa > ptg001221l.fa
8 |
9 | #Isolating the mitochondria
10 | seqkit subseq ptg001221l.fa -r 563637:579472 > cput_mit.fa
11 |
12 | #Removing the mitochondria and breaking the scaffold in two
13 | seqkit subseq ptg001221l.fa -r 1:553681 > cput_ptg1221l_1.fa
14 | seqkit subseq ptg001221l.fa -r 583500:585720 > cput_ptg1221l_2.fa
15 |
16 | #Renaming scaffold ids
17 | sed -i 's/>ptg001221l/>ptg001221l_1/' cput_ptg1221l_1.fa
18 | sed -i 's/>ptg001221l/>ptg001221l_2/' cput_ptg1221l_2.fa
19 | sed -i 's/>ptg001221l/>cput_mit/' cput_mit.fa
20 |
21 | # Joining scaffolds to genome again
22 | cat cput_N_genome.fa cput_ptg1221l_1.fa cput_ptg1221l_2.fa > cput_N_genome_final.fa
23 |
24 | #Checking
25 | grep ">ptg001221l" cput_N_genome_final.fa
26 | grep -c ">" cput_N_genome_final.fa

```

MITOS Web Server

We ran the mitochondrial genome annotation using MITOS2 web server with all the default parameters but the genetic code, which was specified to be the invertebrate one.

We then downloaded the output files to a local computer and sent them to 2-MITOS_results (/home/blowflies/genome_annotation/cput/3-mitochondrial_genome/2-MITOS_results).

The job settings were:

Job ID: cput

Property	Value
Reference	RefSeq 63 Metazoa
Genetic Code	5
Proteins	True
tRNAs	True
rRNAs	True
OH	True
OL	True
Circular	True
Use Al Arab et al.	False
E-value Exponent	2.0
Final Maximum Overlap	50nt
Fragment Quality Factor	100.0
Standard Code	False
Cutoff	50.0%
Clipping Factor	10.0
Fragment Overlap	20.0%
Local only	True
Sensitive only	False
ncRNA overlap:	50 nt

TEST on pacbio genome before assembly

```
1 /home/blowflies/genome_annotation/cput/3-mitochondrial_genome
2 mkdir 2-blast_pacbio
3 # we copied the fastq with the raw pacbio sequences to 2-blast_pacbio
4 gzip -d cput_raw_genome.fastq.gz
5 blastn -task blastn
6 #converting fastq to fasta
7 seqkit fq2fa cmeg_raw_genome.fastq -o cmeg_raw_genome.fa
8 #blastn
9 blastn -task blastn -evalue 0.00001 -db ../1-blastn/cput_mit_database -query
```

4- RepeatModeler

Installing with docker

It was complicated installing all the programs, so we used a Docker container

```
1 #Always use this before repeat modeler
2 docker run -it --rm dfam/tetools:latest
3
4 container-ID: 918a311e45cb
5 container-name: goofy_buck
6 container-image: dfam/tetools:latest
7
8 #To attach the container and continue running press CTRL+P, then CTRL+Q
```

Creating database

```
1 # Moving the fasta file to docker container from the Rosalind server using:
2 # docker cp file.txt container-name:/path/to/copy/file.txt
3 docker cp ./cput_N_genome_final.fa goofy_buck:/home
4
5 # Getting inside the container
6 docker exec -it goofy_buck /bin/bash
7
8 #Database
9 BuildDatabase -name cput_database cput_N_genome_final.fa
```

Running

```
1 RepeatModeler -database cput_database -threads 20 -LTRstruct >log 2>err
2
3 #Then to exit the container and continue running press CTRL+C
```

Copying the results

```
1 # Compressing files inside docker
2 tar -cjvf cput_other_files.tar.gz *
3
4 # now here
5 cd /home/cunha/01-RepeatModeler/cput
6
7 docker cp -a goofy_buck:/home/cput_database-families.stk ./
8 docker cp -a goofy_buck:/home/cput_database-families.fa ./
9 docker cp -a goofy_buck:/home/cput_database-rmod.log ./
10 docker cp -a goofy_buck:/home/log ./
11 docker cp -a goofy_buck:/home/err ./
12 docker cp -a goofy_buck:/home/cput_other_files.tar.gz ./
```

Test - LTR Harvest

First, we split the genome, to see if it works better


```

1 cd /home/blowflies/genome_annotation/cput/0-genome
2
3 seqkit split2 -p 5 cput_N_genome_final.fa
4
5 cd cput_N_genome_final.fa.split/
6
7 for i in *; do mv $i cput_${i##cput_N_genome_final.part_00}; done

```

Now, running LTRHarvest

```

1 # Doing it on Rosalind
2 # Container: cranky_haibt
3
4 cd /home/blowflies/genome_annotation/cput/0-genome/cput_N_genome_final.fa.split/
5
6 for i in *; do docker cp $i cranky_haibt:/home; done
7
8 docker exec -it cranky_haibt /bin/bash
9
10 cd home/
11 mkdir 1 2 3 4 5
12 mv *1.fa 1/
13 mv *2.fa 2/
14 mv *3.fa 3/
15 mv *4.fa 4/
16 mv *5.fa 5/
17
18 # LTR.sh
19 cd /home/1
20 /opt/genometools/bin/gt suffixerator -db cput_1.fa -indexname cput_1_db -tis
21 /opt/genometools/bin/gt ltrharvest -index cput_1_db -minlenltr 100 -maxlenltr
22 /opt/LTR_retriever/LTR_retriever -genome cput_1.fa -inharvest cput_1.harvest.
23
24 cd /home/2
25 /opt/genometools/bin/gt suffixerator -db cput_2.fa -indexname cput_2_db -tis
26 /opt/genometools/bin/gt ltrharvest -index cput_2_db -minlenltr 100 -maxlenltr
27 /opt/LTR_retriever/LTR_retriever -genome cput_2.fa -inharvest cput_2.harvest.
28
29 cd /home/3
30 /opt/genometools/bin/gt suffixerator -db cput_3.fa -indexname cput_3_db -tis
31 /opt/genometools/bin/gt ltrharvest -index cput_3_db -minlenltr 100 -maxlenltr
32 /opt/LTR_retriever/LTR_retriever -genome cput_3.fa -inharvest cput_3.harvest.
33
34 cd /home/4
35 /opt/genometools/bin/gt suffixerator -db cput_4.fa -indexname cput_4_db -tis
36 /opt/genometools/bin/gt ltrharvest -index cput_4_db -minlenltr 100 -maxlenltr
37 /opt/LTR_retriever/LTR_retriever -genome cput_4.fa -inharvest cput_4.harvest.
38
39 cd /home/5
40 /opt/genometools/bin/gt suffixerator -db cput_5.fa -indexname cput_5_db -tis
41 /opt/genometools/bin/gt ltrharvest -index cput_5_db -minlenltr 100 -maxlenltr
42 /opt/LTR_retriever/LTR_retriever -genome cput_5.fa -inharvest cput_5.harvest.
43
44 cd /home

```

It worked! Now we need to download the final files back from the container, concatenate them with the output from the previous RepeatModeler run, and then use this final file to mask the genome

```
1 cd /home/blowflies/genome_annotation/cput/4-RepeatModeler
2 # cput_database-families.fa is here
3
4 docker cp -a cranky_haibt:/home/1/cput_1.fa.LTRlib.fa ./
5 docker cp -a cranky_haibt:/home/2/cput_2.fa.LTRlib.fa ./
6 docker cp -a cranky_haibt:/home/3/cput_3.fa.LTRlib.fa ./
7 docker cp -a cranky_haibt:/home/4/cput_4.fa.LTRlib.fa ./
8 docker cp -a cranky_haibt:/home/5/cput_5.fa.LTRlib.fa ./
9
10 docker cp -a cranky_haibt:/home/cput_LTR.tar.gz ./ # other outputs
11
12 docker stop cranky_haibt # kill container
13
14 cat *fa > cput_modeler_complete.fa
15
16 # remove redundancy
17 vsearch --cluster_fast cput_modeler_complete.fa -id 0.80 -threads 15 -centro
```

5- RepeatMasker

Installing steps can be found in the *Chrysomya megacephala* log (https://hackmd.io/ziOztK1MQZecVm0_qDAr1g)

##Doing it on Rosalind

```
1 # from Darwin to Rosalind
2 cd /home/cunha/01-RepeatModeler/cput
3
4 scp -P 2205 * pedro@143.107.244.181:/home/blowflies/genome_annotation/cput/4-
5
6 # running
7 cd /home/blowflies/genome_annotation/cput/5-RepeatMasker
8
9 RepeatMasker -lib /home/blowflies/genome_annotation/cput/4-RepeatModeler/cput
```

2.5.1 BUSCO

```
1 cd /home/blowflies/genome_annotation/cput/5-RepeatMasker
2 mkdir 1-BUSCO
3 cp cput_N_genome_final.fa.masked 1-BUSCO/
4 cd 1-BUSCO
5 sudo docker run -u $(id -u) -v $(pwd):/busco_wd ezlabgva/busco:v5.4.7_cv1 bus
```

#Results:

```
-----  
|Results from dataset diptera_odb10|  
-----  
|C:98.8%[S:93.6%,D:5.2%],F:0.2%,M:1.0%,n:3285|  
|3246   Complete BUSCOs (C)|  
|3075   Complete and single-copy BUSCOs (S)|  
|171    Complete and duplicated BUSCOs (D)|  
|6       Fragmented BUSCOs (F)|  
|33      Missing BUSCOs (M)|  
|3285   Total BUSCO groups searched|  
-----
```

2.5.2 QUAST

```
1 | #comparing with the file /home/blowflies/genome_annotation/cput/0-genome/cput  
2 | cd /home/blowflies/genome_annotation/cmeg/5-RepeatMasker  
3 | mkdir 2-QUAST  
4 | cp cput_N_genome_final.fa.masked 2-QUAST/  
5 | cd 2-QUAST  
6 | mkdir cput.fa  
7 | cd cput.fa  
8 | sudo quast.py ../cput_N_genome_final.fa.masked -t 20 --eukaryote --large --rr
```

```
1 | #Results:  
2 | cd /home/blowflies/genome_annotation/cput/5-RepeatMasker/2-QUAST/cput.fa/quas  
3 | cat report.txt
```

```

Assembly                                cput_N_genome_final.fa.masked
# contigs (>= 0 bp)                     653
# contigs (>= 1000 bp)                   653
# contigs (>= 5000 bp)                   652
# contigs (>= 10000 bp)                  650
# contigs (>= 25000 bp)                  551
# contigs (>= 50000 bp)                  446
Total length (>= 0 bp)                   587036998
Total length (>= 1000 bp)                 587036998
Total length (>= 5000 bp)                 587034777
Total length (>= 10000 bp)                587016365
Total length (>= 25000 bp)                585212355
Total length (>= 50000 bp)                581400726
# contigs                                652
Largest contig                           15138138
Total length                             587034777
Estimated reference length                5000000000
GC (%)                                    28.98
N50                                        4322519
NG50                                       4997734
N75                                        1939160
NG75                                       2854789
L50                                        45
LG50                                       35
L75                                        95
LG75                                       68
# total reads                            768
# left                                    0
# right                                   0
Mapped (%)                               100.0
Properly paired (%)                      0.0
Avg. coverage depth                      1
Coverage >= 1x (%)                      100.0
# N's per 100 kbp                       0.07
# predicted rRNA genes                   238 + 31 part

```

6- RNA-seq

We extracted RNA from:

- 50 eggs
- 10 L1
- 5 L2
- 2 L3
- 1 pupae
- 1 virgin female
- 1 gravid female
- 1 male

Then, we pooled all the samples (2ug of RNA from each sample) and sequenced it.

- RNAseq Illumina 20M reads paired end PE150 Q30>85%

```

1  #Coping files
2  cp -r /home/Raw_seqs/cput_pool_RNA /home/blowflies/genome_annotation/cput/
3  #Renaming
4  cd /home/blowflies/genome_annotation/cput/
5  mv cput_pool_RNA 8-cput_pool_RNA
6  #checking md5
7  cd 8-cput_pool_RNA/
8  cat MD5.txt
9  #8ecb8a84268de8259d50a885c316d6a9  Cput_1.fq.gz
10 #d6c794f0578c674b16f3f65c5da3933b  Cput_2.fq.gz
11 md5sum Cput*
12 #8ecb8a84268de8259d50a885c316d6a9  Cput_1.fq.gz
13 #d6c794f0578c674b16f3f65c5da3933b  Cput_2.fq.gz
14 mkdir 0-raw_reads
15 mv Cput* 0-raw_reads
16 mv MD5.txt 0-raw_reads

```

Quality control of raw reads

We ran FastQC and, then MultiQC.

Don't need to unzip raw read files because fastqc can cope with zipped files (.gz).

FastQC will process one sample at a time and give you an output report for each sample separately. MultiQC will combine all the outputs from FastQC analysis and give you one QC report for all processed samples, making them more easily comparable.

→ nice webpage on fastqc and multiqc: https://stab.st-andrews.ac.uk/wiki/index.php/FASTQC_and_MultiQC

→ <https://multiqc.info/>

```

cd /home/blowflies/genome_annotation/cput/8-cput_pool_RNA
mkdir 1-QC
cd 0-raw_reads
fastqc *fq.gz #v0.11.9
multiqc . #Version 1.11
mv *.html ../1-QC
mv *.zip ../1-QC
mv multiqc_data ../1-QC

```

Coping multiqc report to a local computer

```

scp -P 2205 vanessa@143.107.244.181:/home/blowflies/genome_annotation/cput/8-
cput_pool_RNA/1-QC/multiqc_report.html /mnt/c/Users/vansc/Downloads

```

Results:

https://drive.google.com/file/d/1s44TfwlPpu2_-2_vLqoiHMEY9jyUq3-l/view?usp=share_link

Trimming

Processing raw reads to trimming (remove only bad quality bases).

I used Trimmomatic to trimming version 0.39

→ nice webpage on how to use Trimmomatics: <http://www.usadellab.org/cms/index.php?page=trimmomatic>

<https://datacarpentry.org/wrangling-genomics/03-trimming/>

```
cd /home/blowflies/genome_annotation/cput/8-cput_pool_RNA/0-raw-reads
mkdir ../2-trimming

TrimmomaticPE Cput_1.fq.gz Cput_2.fq.gz -threads 8 -baseout
/home/blowflies/genome_annotation/cput/8-cput_pool_RNA/2-
trimming/cput.trimmed.fq.gz ILLUMINACLIP:TruSeq3-PE-2.fa:2:30:10
SLIDINGWINDOW:4:15 MINLEN:36
```

Quality control of trimmed reads

```
# in this directory -> /home/blowflies/genome_annotation/cput/4-cput_pool_RNA/2-tr
# QC
fastqc *.gz

# here -> /home/blowflies/genome_annotation/calb/4-calb_pool_RNA/1-QC
mkdir trimmed_reads_qc

# in this directory -> /home/blowflies/genome_annotation/cput/4-cput_pool_RNA/2-tr
mv *.html ../1-QC/trimmed_reads_qc/
mv *.zip ../1-QC/trimmed_reads_qc/
multiqc .
```

Observation

We renamed the genome file from this step onwards.

```
1 | cd /home/blowflies/genome_annotation/cput/0-genome
2 |
3 | mv purged.fa cput.fa
```

6.1 - New RNA sequencing

As we did not have a great representation of BUSCO genes and low quality in general in our previous sequencing, we did it again. The new data can be found at:

```
1 | /home/blowflies/genome_annotation/cput/8-cput_pool_RNA/new_seq
```

Quality control of raw reads

We ran FastQC and, then MultiQC.

Don't need to unzip raw read files because fastqc can cope with zipped files (.gz).

FastQC will process one sample at a time and give you an output report for each sample separately. MultiQC will combine all the outputs from FastQC analysis and give you one QC report for all processed samples, making them more easily comparable.

→ nice webpage on fastqc and multiqc: https://stab.st-andrews.ac.uk/wiki/index.php/FASTQC_and_MultiQC

→ <https://multiqc.info/>

```
cd /home/blowflies/genome_annotation/cput/8-cput_pool_RNA/new_seq
mkdir 1-QC
cd 0-raw_reads
fastqc *fastq.gz #v0.11.9
multiqc . #Version 1.11
mv *.html ../1-QC
mv *.zip ../1-QC
mv multiqc_data ../1-QC
```

Coping multiqc report to a local computer

```
scp -P 2205 diniz@143.107.244.181:/home/blowflies/genome_annotation/cput/8-cput_pool_RNA/new_seq/1-QC/multiqc_report.html /Users/diniz/Desktop/
```

Results:

https://drive.google.com/drive/folders/1g_VW5o6HiXAKChcWKe_GDxqSjqAMdGAn

Trimming

Processing raw reads to trimming (remove only bad quality bases).

I used Trimmomatic to trimming version 0.39

→ nice webpage on how to use Trimmomatics: <http://www.usadellab.org/cms/index.php?page=trimmomatic>

<https://datacarpentry.org/wrangling-genomics/03-trimming/>

```
cd /home/blowflies/genome_annotation/cput/8-cput_pool_RNA/0-raw-reads
mkdir ../2-trimming

TrimmomaticPE Cupt1_NGS629_S31_L001_R1_001.fastq.gz
Cupt1_NGS629_S31_L001_R2_001.fastq.gz -threads 8 -baseout
/home/blowflies/genome_annotation/cput/8-cput_pool_RNA/new_seq/2-
```

```
trimming/cput.trimmed.fq.gz ILLUMINACLIP:TruSeq3-PE-2.fa:2:30:10
SLIDINGWINDOW:4:15 MINLEN:36
```

Quality control of trimmed reads

```
# in this directory -> /home/blowflies/genome_annotation/cput/4-cput_pool_RNA/2-tr
# QC
fastqc *.gz

# here -> /home/blowflies/genome_annotation/calb/4-calb_pool_RNA/1-QC
mkdir trimmed_reads_qc

# in this directory -> /home/blowflies/genome_annotation/cput/4-cput_pool_RNA/2-tr
mv *.html ../1-QC/trimmed_reads_qc/
mv *.zip ../1-QC/trimmed_reads_qc/
mv multiqc_data ../1-QC/trimmed_reads_qc/
multiqc .
```

Assembling the transcriptome ("old + new" RNA-seq)

In /home/cunha/03-RNA/01-Reads/cput/new_seq

```
1  #!/bin/bash
2
3  #SBATCH --job-name trinity ## nome que aparecerá na fila
4  #SBATCH --output trinity_calb.out ## nome do arquivo de saída; o %j é igual a
5  #SBATCH --ntasks=1 ## número de tarefas (análises) a serem executadas
6  #SBATCH --cpus-per-task=20 ## o número de threads alocados para cada tarefa
7  #SBATCH --mem-per-cpu=1000M # memória por núcleo da CPU
8  #SBATCH --partition=long ## as partições a serem executadas (separadas por v
9  #SBATCH --time=10-00:00:00 ## hora para análise (dia-hora:min:seg)
10 #SBATCH --error=err
11
12 srun docker run --rm -v`pwd`:`pwd` trinityrnaseq/trinityrnaseq Trinity --seq
```

BUSCO of the new transcriptome

In rosalind (Darwin was not working) /home/diniz/cput_all_busco

```
1  docker run -u $(id -u) -v $(pwd):/busco_wd ezlabgva/busco:v5.4.7_cv1 busco -i
```

7- Trinity

Transcriptome assembly

Moving trimmed reads to Darwin


```

1 | mkdir /home/cunha/03-RNA/01-Reads/cmeg
2 |
3 | scp -P 4988 /home/blowflies/genome_annotation/cput/8-cput_pool_RNA/2-trimming

```

Assembling the transcriptome

```

1 | cd /home/cunha/03-RNA/01-Reads/cput
2 |
3 | docker run --rm -v `pwd`:`pwd` trinityrnaseq/trinityrnaseq Trinity --seqType
4 |
5 | mv cput_trinity.* /home/cunha/03-RNA/02-Trinity/cput_trinity
6 |
7 | # Transcriptome size
8 | grep -c ">" cput_trinity.Trinity.fasta # 28967

```

BUSCO (transcriptome quality)

```

1 | mkdir /home/cunha/03-RNA/03-Busco
2 |
3 | docker pull ezlabgva/busco:v5.4.4_cv1 # just because we didn't have busco on
4 |
5 | # we copied all transcriptomes in this directory and ran everything at once
6 | for i in *; do docker run -u $(id -u) -v $(pwd):/busco_wd ezlabgva/busco:v5.4

```

Transcriptome	Complete (all)	Complete Single	Complete Dup.	Fragmented	Missing
only old	24.9	17.7	7.2	9.1	66
only new	67.7	50	17.7	8.4	23.9
old + new	69.4	45.1	24.3	8.5	22.1

8- STAR - RNAseq alignment

index

All genome FASTA files **cannot** be zipped

```

1 # unzipping files
2 cd
3 gzip -dk *P*
4
5 # We need to create a directory where the genome indexes will be stored before
6 cd /home/blowflies/genome_annotation/cput
7 mkdir 9-STAR
8 chmod 777 9-STAR
9 cd 9-STAR
10 mkdir star_index
11 chmod 777 star_index
12
13 STAR --runThreadN 20 --runMode genomeGenerate --genomeDir /home/blowflies/gen

```

Before the alignment itself we had to concatenate all the fastq files available. The files are in Rosalind (/home/blowflies/genome_annotation/cput/8-cput_pool_RNA/2-trimming/all)

```

1 cat *_1P.fq.gz > cput_all_1P.fq.gz
2 cat *_2P.fq.gz > cput_all_2P.fq.gz

```

STAR alignment

```

1 cd /home/blowflies/genome_annotation/cput/8-cput_pool_RNA/all
2
3 for i in *_1P.fq.gz; do STAR --runMode alignReads --readFilesCommand zcat --c

```

9- Braker3

All information is in the Calb file (https://hackmd.io/KMzfBC2aQ9qliTy11f_QhQ)

Final outputs are here: /home/blowflies/genome_annotation/cput/10-BRAKER3

BRAKER3 second run

```

1 # need to run first
2 export BRAKER_SIF=/home/diniz/programs/braker/braker3.sif
3
4 #We copied the file 10-BRAKER3 to the braker file
5 cd /home/diniz/programs/braker/
6 mkdir Cput_2
7 cp /home/blowflies/genome_annotation/cput/10-BRAKER3 /home/diniz/programs/braker/
8 cd /home/diniz/programs/braker/Cput_2/10-BRAKER3
9 mkdir restart
10 # Run
11 singularity exec /home/diniz/programs/braker/braker3.sif braker.pl --genome=,

```

Total transcripts: 21313

Busco evaluation with protein sequences from restart Braker run

```
1 cd /home/diniz/programs/braker/Cput_2/10-BRAKER3/restart
2 mv braker.gtf cput_braker.gtf
3 mv brker.aa cput_braker.aa
4
5 docker run -u $(id -u) -v $(pwd):/busco_wd ezlabgva/busco:v5.4.4_cv1 busco -i
```

```
# BUSCO version is: 5.4.4
# The lineage dataset is: diptera_odb10 (Creation date: 2020-08-05, number of genes: 1000)
# Summarized benchmarking in BUSCO notation for file /busco_wd/cput_braker.aa
# BUSCO was run in mode: proteins
```

```
***** Results: *****
```

```
C:97.7%[S:83.7%,D:14.0%],F:0.9%,M:1.4%,n:3285
3211   Complete BUSCOs (C)
2750   Complete and single-copy BUSCOs (S)
461    Complete and duplicated BUSCOs (D)
29     Fragmented BUSCOs (F)
45     Missing BUSCOs (M)
3285   Total BUSCO groups searched
```

11- EnTAP

We did it on Darwin using the output from the restart Braker run

```
1 #We copied the aminoacid file from Rosalind (/home/diniz/programs/braker/Cput_2/10-BRAKER3/restart)
2 cd /home/martins/EnTAP_restart
3
4 EnTAP --runP -i /home/martins/EnTAP_restart/Proteomes/cput_braker.aa -d /home/martins/EnTAP_restart/Proteomes/
5
6 mv entap_outfiles/ cput/
7
8 scp -r -P 2205 entap_outfiles/ pedro@143.107.244.181:/home/blowflies/genome_annotation/
```

Checking md5

```
1 cd /home/blowflies/genome_annotation/cput/11-EnTAP/entap_outfiles/final_results
2 md5sum entap_results.tsv
3 #ded21bc806abf83077034a891b2e5cdd  entap_results.tsv
4 cd /home/diniz/programs/braker/Cput_2/10-BRAKER3/restart
5 md5sum cput_braker.gtf
6 #ec60e5bfed4626446c95fba0d56a4bd3  cput_braker.gtf
```

Making an unique gtf file with augustus and ENTAP outputs files:

- /home/blowflies/genome_annotation/cput/11-EnTAP/entap_outfiles/final_results/entap_results.tsv
#ded21bc806abf83077034a891b2e5cdd
- /home/diniz/programs/braker/Cput_2/10-BRAKER3/restart/cput_braker.gtf
#ec60e5bfed4626446c95fba0d56a4bd3

In R:

```

1  # matching AUGUSTUS and ENTAP output into a unique gtf
2
3  # libraries
4  library(data.table)
5  library(dplyr)
6
7  # reading the files
8  tsv <- fread(file = "entap_results.tsv", header = FALSE)
9  tsv <- tsv[-1,]
10 tsv <- tsv[,c(1,13)]
11 gtf <- fread(file = "cput_final.gtf")
12
13 # updated gtf
14 new_gtf <- left_join(gtf, tsv, by = c("V9" = "V1"))
15
16 # write gtf
17 fwrite(x = new_gtf, quote = FALSE, sep = '\t', row.names = FALSE,
18        col.names = FALSE, file = "cput_entap_final.gtf")
19
20 # to know how many annotated transcripts are (annotated proteins)
21 ann_tra <- na.omit(tsv$V13) # look the number of elements in this and compare

```

Total annotated transcripts: 16843 (out of 21313)

We copied the final gtf file to Rosalind server and checked md5:

```

1  cd /Users/diniz/Desktop
2  md5 Cput_annot.gtf
3  #0386ed7e643c421d139327e33a81f699
4  scp -P 2205 Cput_annot.gtf diniz@143.107.244.181:/home/blowflies/genome_annotation/
5
6  #On Rosalind:
7  cd /home/blowflies/genome_annotation/cput/12-final_files/
8  md5sum Cput_annot.gtf
9  #74d1c8756206130b81c9db811ad7713a

```

12-Final files

```

1 cd /home/blowflies/genome_annotation/cput/12-final_files
2 md5sum *
3 #0386ed7e643c421d139327e33a81f699 Cput_annot.gtf
4 #87fd55c764a84b7229abf3e493a8c0e5 Cput_cds.fa
5 #2213621d218d62c90311483ddedc2ceb Cput_genome.fa
6 #67d462183e79c464cb5305e593e5f669 Cput_protein.aa
7 ``
8
9 | type |original/copy| file | Path | md5 |
10 | --- | --- | --- | --- | --- |
11 |Raw genome| original| XDOVE_20221110_S64411e_PL100270437-1_D01.ccs.fastq.gz
12 |Raw genome| original| purged.fa| /home/Reference_genomes/Cputoria| 314250c4
13 |Raw genome| copy| cput.fa|/home/blowflies/genome_annotation/cput/0-genome|
14
15 | type |original/copy| file | Path | md5 |
16 | --- | --- | --- | --- | --- |
17 |Mitochondrial genome|original| cput_mit.fa| /home/blowflies/genome_annotation|
18 | Nuclear genome| original| cput_N_genome_final.fa| /home/blowflies/genome_annotation|
19
20 | type |original/copy| file | Path | md5 |
21 | --- | --- | --- | --- | --- |
22 |RNA-seq| original|Cput_1.fq.gz| /home/Raw_seqs/cput_pool_RNA| 8ecb8a84268d
23 | RNA-seq| original|Cput_2.fq.gz| /home/Raw_seqs/cput_pool_RNA| d6c794f0578c
24 |RNA-seq unzipped unzipped| original| Cput_1.fq| /home/Raw_seqs/cput_pool_RNA|
25 |RNA-seq unzipped| original| Cput_2.fq |/home/Raw_seqs/cput_pool_RNA| 8c13e3
26 |Trimmed reads zipped| original| cput_all_1P.fq.gz| /home/blowflies/genome_annotation|
27 |Trimmed reads zipped| original| cput_all_2P.fq.gz| /home/blowflies/genome_annotation|
28
29 ==**Não sei qual é o arquivo correto**==
30 ==|Transcriptome| original| ???| /home/cunha/03-RNA/02-Trinity| ???|==
31
32 | type |original/copy| file | Path | md5 |
33 | --- | --- | --- | --- | --- |
34 | Masked genome | original|cput_N_genome_final.fa.masked | /home/blowflies/genome_annotation|
35 | Masked genome |copy| Cput_masked.fasta | /home/blowflies/genome_annotation|
36 | Masked genome |copy|Cput_masked.fasta| /home/blowflies/genome_annotation/c|
37 | Masked genome |copy|Cput_masked.fasta| /home/diniz/programs/braker/Cput_2/|
38 | Masked genome | copy| Cput_masked.fasta | /home/pedro/Non\_Coding\_Element|
39
40 | type |original/copy| file | Path | md5 |
41 | --- | --- | --- | --- | --- |
42 |Proteome |original| cput_braker.aa | /home/diniz/programs/braker/Cput_2/10-|
43 |Proteome |copy| cput_braker.aa | /home/martins/EntAP_restart/Proteomes | 67|
44 |Proteome |copy| cput_braker.aa | /home/00-Sequences/Cputoria/01-Genomic_data|
45
46 | type |original/copy| file | Path | md5 |
47 | --- | --- | --- | --- | --- |
48 |Codingseq| original| braker.codingseq | /home/diniz/programs/braker/Cput_2|
49 |Codingseq|copy| cput_braker.codingseq | /home/00-Sequences/Cputoria/01-Genomic_data|
50
51 | type |original/copy| file | Path | md5 |
52 | --- | --- | --- | --- | --- |
53 |gtf output Braker3 second run| original| cput_braker.gtf| /home/diniz/programs/braker/Cput_2|
54 |gtf output Braker3 second run| copy| | **Computador do Diniz** | |
55 |Final gtf| original| | **Computador do Diniz** | |
56 |Final gtf| copy| | **Rosalind** | |
57
58 **FALTA: conferir md5 dos arquivos usados no script do R do entap e cmeg_entap**
59

```

```

60  **PROBLEMAS:**
61  - **arquivo proteoma da pasta 00-sequences está com md5 diferente do original**
62  - **arquivo codingseq da pasta 00-sequences está com md5 diferente do original**
63  - **arquivo gtf que está na pasta 00-sequences não é o final**
64
65  -----
66  # OBSOLETE
67  # 10- gFACs
68  https://gfacs.readthedocs.io/en/latest/Flags/index.html
69
70  ```bash=
71  cd /home/blowflies/genome_annotation/cput/
72  mkdir 11-gfacs
73  cd /home/blowflies/genome_annotation/cput/11-gfacs
74  mkdir results
75  sudo perl /gFACs-master/gFACs.pl -f braker_2.1.2_gtf -p cput --rem-all-incon
76  ```
77
78  Results:
79  Number of genes (Augustus/BRAKER): 21450
80  Number of genes (gFACs): 20894
81
82  # 11- EnTAP
83  We did it on Darwin
84  ```bash=
85  cd /home/martins/EnTAP
86
87  # The gFACs outputs for all species are here (*_genes.fasta.faa)
88  mkdir Proteomes
89
90  mkdir cput
91  cd cput
92  EnTAP --runP -i /home/martins/EnTAP/Proteomes/cput_genes.fasta.faa -d /home/
93  ```
94  # 12. gFACs again (with EnTAP output)
95  ```bash=
96  cd /home/blowflies/genome_annotation/cput/
97  mkdir 13-gfacs_entap
98  cd /home/blowflies/genome_annotation/cput/13-gfacs_entap
99  mkdir results
100 sudo perl /gFACs-master/gFACs.pl -f gFACs_gene_table -p cput --rem-all-incon
101 ```
102
103 # 13. Final annotation
104 Final files are here:
105 ```bash=
106 /dados/home/blowflies/genome_annotation/cput/14-final_annot
107 # 24543 gene models
108 ```
109
110 Final busco
111 ```bash=
112 docker run -u $(id -u) -v $(pwd):/busco_wd ezlabgva/busco:v5.4.4_cv1 busco -
113 ```
114 
115
116 # OBSOLETE STUFF
117 ## RepeatMasker
118

```

```

119  ````bash=
120  # we start here
121  /home/blowflies/genome_annotation/cput
122
123  # new directory
124  mkdir 7-RepeatMasker
125  cd 7-RepeatMasker/
126  mkdir With_drosophila
127  cd With_drosophila/
128
129  # running
130  RepeatMasker -species "Drosophila melanogaster" -dir . -pa 8 -a -xsmall -s

```

14 - Submission to NCBI

NCBI require a sqn file for assembly submission. I followed the step described in

https://www.ncbi.nlm.nih.gov/genbank/genomes_gff/

Then I had to rename the headers in the genome fasta and convert from gtf to gff

```

1  sed -i '/^>/ s/$/ [organism=Chrysomya putoria]/' Cput_genome.fa
2  singularity run agat_1.0.0--pl5321hdfd78af_0.sif
3  agat_convert_sp_gxf2gxf.pl --gtf Cput_annot.gtf --output Cput_annot.gff

```

Finally I ran table2asn to get the sqn file

```

1  /home/diniz/programs/linux64.table2asn -M n \
2  -J \
3  -c w \
4  -euk \
5  -gaps-min 10 \
6  -f Cput_annot.gff \
7  -i Cput_genome.fa \
8  -locus-tag-prefix Cput \
9  -o cput.sqn \
10 -Z \
11 -V b

```