

# Log: Genome Annotation Chrysomya megacephala

---

tags: [Genome Annotation](#) [Cmeg](#)

## Table of Contents

---

- [Log: Genome Annotation Chrysomya megacephala](#)
  - [Table of Contents](#)
  - [Genome annotation workflow:](#)
  - [List of software:](#)
- [Genome sequencing](#)
- [1- BUSCO First](#)
- [2- QUAST](#)
  - [Running](#)
  - [View](#)
- [3- Mitochondrial genome](#)
  - [Blastn](#)
  - [MITOS Web Server](#)
  - [TEST on pacbio genome before assembly](#)
  - [Pairwise Sequence Alignment Tools \(Reverse complement + Trimming the mt genome\)](#)
- [4- RepeatModeler](#)
  - [Installing with docker](#)
  - [Creating database](#)
  - [Running](#)
  - [Copying the results](#)
  - [Test - LTR Harvest](#)
- [5- RepeatMasker](#)
  - [5.1 BUSCO](#)
  - [5.2 QUAST](#)
- [6- RNA-seq](#)
  - [Quality control of raw reads](#)
  - [Trimming](#)
    - [Quality control of trimmed reads](#)
    - [Observation](#)
    - [Quality Control of raw reads](#)
    - [Trimming](#)
- [7 - Trinity](#)
  - [Running Trinity](#)
  - [Assembly statistics](#)
  - [BUSCO \(transcriptome quality\)](#)
- [8- STAR - RNAseq alignment](#)
  - [Index](#)
  - [STAR alignment](#)
- [9- BRAKER3 - Structural annotation](#)
  - [BRAKER3 run](#)
  - [Busco evaluation with protein sequences](#)
  - [BRAKER3 second run](#)

- Busco evaluation with protein sequences from second Braker run
- 11- EnTAP
- 12- Final files
- OBSOLETE
- 10- gFACs
- 11- EnTAP
- 12. gFACs again (with EnTAP output)
- 13. Final annotation
- OBSOLETE STUFF
  - RepeatModeler
  - RNA reads from NCBI
- Trinity (first)
  - Moving trimmed reads to Darwin
  - Assembling the transcriptome
  - BUSCO (transcriptome quality)
- 14 - Submission to NCBI

---

## 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/Genome_Assembly)

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

---

## List of software:

- ☒ QUAST (/usr/local/bin/quast.py) (Pedro!)
- ☒ RepeatModeler (/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>

☒ GeneMark-ETP (PEDRO!)

- <https://github.com/gatech-genemark/GeneMark-ETP>

---

## Genome sequencing

---

We sequenced 1 male at Dovetail (Pac-Bio).

### ANOTAR AQUI DADOS DO SEQUENCIAMENTO

**Number of reads** 2521875

**Coverage (x)** 66

**HPA Length (bp)** 1446676375

**HPA N50 (bp)** 1148833

**HPA N90 (bp)** 117041

**HPA L50** 287

**HPA L90** 1800

**FA Length (bp)** 671207201

**FA N50 (bp)** 2214294

**FA N90 (bp)** 540686

**FA L50** 84

**FA L90** 316

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.

---

## 1- BUSCO First

---

Busco version 5.3.2

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

---

## 2- QUAST

---

Quast accepts assemblies and reference genomes in FASTA format. Files may be compressed with zip, gzip or bzip2.

Quast accepts Illumina, PacBio, and Oxford Nanopore reads in FASTQ format (may be compressed).

---

## Running

---

```

1  pwd
2  /home/blowflies/genome_annotation/cmeg/2-quast
3
4  # OLD
5  # /quast-5.0.2/quast.py ../0-genome/purged.fa -t 20 --eukaryote --large --rna-finding
6
7  # Now we're running from the python library
8
9  # Obs: had to edit the jsontemplate.py script again to change cgi to html
10 # We're in the root
11 sudo find / -type d -name "quast_libs"
12
13 # the script we need to edit is here:
14 /usr/local/lib/python3.8/dist-packages/quast-5.0.2-py3.8.egg/quast_libs/site-packages/quast/jsontemplate.py
15
16 # Running
17 sudo quast.py ../0-genome/purged.fa -t 20 --eukaryote --large --rna-finding
18

```

## View

```
less quast_results/latest/report.txt
```

## 3- Mitochondrial genome

Cmeg mitochondrial genome on NCBI:

[https://www.ncbi.nlm.nih.gov/nucleotide/NC\\_019633.1](https://www.ncbi.nlm.nih.gov/nucleotide/NC_019633.1)

NC\_019633.1

```

1  cd home/blowflies/genome_annotation/cmeg/3-mitochondrial_genome
2  mkdir 1-blastn
3  # We downloaded mitochondrial sequence from NCBI to a local computer and tra

```

## Blastn

version 2.11

```

1  cd home/blowflies/genome_annotation/cmeg/3-mitochondrial_genome/1-blastn
2  #Making database using the genome
3  makeblastdb -in ../Cmeg_ref_mitochondria.fasta -dbtype nucl -out cmeg_mit_database
4  #Running
5  blastn -task blastn -evalue 0.00001 -db ./cmeg_mit_database -query ../0-genome

```

The best alignment was against scaffold 268, which is a very big scaffold (1.9Mb). We isolated this scaffold and cut it in three pieces "before mitochondria" - "mitochondria" - "after mitochondria"

```

1  cd /home/blowflies/genome_annotation/cmeg/0-genome
2
3  # Remove scaffold from genome
4  seqkit grep -v -p "ptg0002681" cmeg.fa > cmeg_N_genome.fa
5
6  # Isolate scaffold
7  seqkit grep -p "ptg0002681" cmeg.fa > ptg0002681.fa
8
9  # Isolating the mitochondria
10 seqkit subseq ptg0001681.fa -r 1324186:1339448 > cmeg_mit.fa
11
12 #Removing the mitochondria and breaking the scaffold in two
13 seqkit subseq ptg0002681.fa -r 1:1305525 > cmeg_ptg2681_1.fa
14 seqkit subseq ptg0002681.fa -r 1341184:1921755 > cmeg_ptg2681_2.fa
15
16 #Renaming scaffold ids
17 sed -i 's/>ptg0002681/>ptg0002681_1/' cmeg_ptg2681_1.fa
18 sed -i 's/>ptg0002681/>ptg0002681_2/' cmeg_ptg2681_2.fa
19 sed -i 's/>ptg0002681/>cmeg_mit/' cmeg_mit.fa
20
21 # Joining scaffolds to genome again
22 cat cmeg_N_genome.fa cmeg_ptg2681_1.fa cmeg_ptg2681_2.fa > cmeg_N_genome_final.fa
23
24 #Checking
25 grep ">ptg0002681" cmeg_N_genome_final.fa
26 grep -c ">" cmeg_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/cmeg/3-mitochondrial\_genome/2-MITOS\_results).

The job settings were:

Job ID: cmeg	
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/cmeg/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 cmeg_raw_genome.fastq.gz
5 #converting fastq to fasta
6 seqkit fq2fa cmeg_raw_genome.fastq -o cmeg_raw_genome.fa
7 #blastn
8 blastn -task blastn -evalue 0.00001 -db ../1-blastn/cmeg_mit_database -query
```

## Pairwise Sequence Alignment Tools (Reverse complement + Trimming the mt genome)

```
1 #We are here
2 cd /home/blowflies/genome_annotation/cmeg/0-genome
3
4 #First rough trimming
5 seqkit subseq ptg0002681.fa -r 1305525:1341184 > cmeg_mit_repeats.fa
6
7 #Downloaded to local computer
8 scp -P 2205 diniz@143.107.244.181:/home/blowflies/genome_annotation/cmeg/0-genome
```

<https://www.ebi.ac.uk/Tools/psa/>

- Water (EMBOSS)  
EMBOSS Water uses the Smith-Waterman algorithm (modified for speed enhancements) to calculate the local alignment of two sequences.
- Stretcher (EMBOSS)  
EMBOSS Stretcher uses a modification of the Needleman-Wunsch algorithm that allows larger sequences to be globally aligned.
- Matcher (EMBOSS)  
EMBOSS Matcher identifies local similarities between two sequences using a rigorous algorithm based on the LALIGN application.

## 4- RepeatModeler

Couldn't run without docker!!! I don't know what is happening (Vanessa)

## 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: 1b8109bcf5aa
5 container-name: hungry_banach
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 ./cmeg_N_genome_final.fa hungry_banach:/home/cmeg
4
5 # Getting inside the container
6 docker exec -it hungry_banach /bin/bash
7
8 #Database
9 BuildDatabase -name cmeg_database cmeg_N_genome_final.fa
```

## Running

```
1 RepeatModeler -database cmeg_database -threads 30 -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 cmeg_other_files.tar.gz *
3
4 # now here
5 cd /home/cunha/01-RepeatModeler/cmeg
6
7 docker cp -a hungry_banach:/home/cmeg/cmeg_database-families.stk ./
8 docker cp -a hungry_banach:/home/cmeg/cmeg_database-families.fa ./
9 docker cp -a hungry_banach:/home/cmeg/cmeg_database-rmod.log ./
10 docker cp -a hungry_banach:/home/cmeg/log ./
11 docker cp -a hungry_banach:/home/cmeg/err ./
12 docker cp -a hungry_banach:/home/cmeg/cmeg_other_files.tar.gz ./
```

## Test - LTR Harvest

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

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

file	format	type	num_seqs	sum_len	min_len	avg_len	max_len
cmeg_1.fa	FASTA	DNA	152	118,059,414	11,549	776,706.7	8,638,607
cmeg_2.fa	FASTA	DNA	152	159,414,450	13,769	1,048,779.3	9,257,387
cmeg_3.fa	FASTA	DNA	152	142,440,787	12,578	937,110.4	14,378,473
cmeg_4.fa	FASTA	DNA	152	137,312,730	11,578	903,373.2	7,470,457
cmeg_5.fa	FASTA	DNA	152	113,944,162	11,647	749,632.6	6,742,468

Now, running LTRHarvest

```

1  # Doing it on Rosalind
2  # Container: cranky_morse
3
4  cd /home/blowflies/genome_annotation/cmeg/0-genome/cmeg_N_genome_final.fa.sp
5
6  for i in *; do docker cp $i cranky_morse:/home; done
7
8  docker exec -it cranky_morse /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 cmeg_1.fa -indexname cmeg_1_db -tis
21 /opt/genometools/bin/gt ltrharvest -index cmeg_1_db -minlenltr 100 -maxlenltr
22 /opt/LTR_retriever/LTR_retriever -genome cmeg_1.fa -inharvest cmeg_1.harvest
23
24 cd /home/2
25 /opt/genometools/bin/gt suffixerator -db cmeg_2.fa -indexname cmeg_2_db -tis
26 /opt/genometools/bin/gt ltrharvest -index cmeg_2_db -minlenltr 100 -maxlenltr
27 /opt/LTR_retriever/LTR_retriever -genome cmeg_2.fa -inharvest cmeg_2.harvest
28
29 cd /home/3
30 /opt/genometools/bin/gt suffixerator -db cmeg_3.fa -indexname cmeg_3_db -tis
31 /opt/genometools/bin/gt ltrharvest -index cmeg_3_db -minlenltr 100 -maxlenltr
32 /opt/LTR_retriever/LTR_retriever -genome cmeg_3.fa -inharvest cmeg_3.harvest
33
34 cd /home/4
35 /opt/genometools/bin/gt suffixerator -db cmeg_4.fa -indexname cmeg_4_db -tis
36 /opt/genometools/bin/gt ltrharvest -index cmeg_4_db -minlenltr 100 -maxlenltr
37 /opt/LTR_retriever/LTR_retriever -genome cmeg_4.fa -inharvest cmeg_4.harvest
38
39 cd /home/5
40 /opt/genometools/bin/gt suffixerator -db cmeg_5.fa -indexname cmeg_5_db -tis
41 /opt/genometools/bin/gt ltrharvest -index cmeg_5_db -minlenltr 100 -maxlenltr
42 /opt/LTR_retriever/LTR_retriever -genome cmeg_5.fa -inharvest cmeg_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/cmeg/4-RepeatModeler
2  # cmeg_database-families.fa is here
3
4  docker cp -a cranky_morse:/home/1/cmeg_1.fa.LTRlib.fa ./
5  docker cp -a cranky_morse:/home/2/cmeg_2.fa.LTRlib.fa ./
6  docker cp -a cranky_morse:/home/3/cmeg_3.fa.LTRlib.fa ./
7  docker cp -a cranky_morse:/home/4/cmeg_4.fa.LTRlib.fa ./
8  docker cp -a cranky_morse:/home/5/cmeg_5.fa.LTRlib.fa ./
9
10 docker cp -a cranky_morse:/home/cmeg_LTR.tar.gz ./ # other outputs
11
12 cat *fa > cmeg_modeler_complete.fa
13
14 # remove redundancy
15 vsearch --cluster_fast cmeg_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](https://hackmd.io/ziOztK1MQZecVm0_qDAr1g))

##Doing it on Rosalind

```

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

```

## 5.1 BUSCO

```

1 cd /home/blowflies/genome_annotation/cmeg/5-RepeatMasker
2 mkdir 1-BUSCO
3 cp cmeg_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 busco

```

### Results

```

=====
|Results from dataset diptera_odb10|
=====
|C:97.8%[S:86.6%,D:11.2%],F:0.2%,M:2.0%,n:3285|
|3214 Complete BUSCOs (C)|
|2846 Complete and single-copy BUSCOs (S)|
|368 Complete and duplicated BUSCOs (D)|
|8 Fragmented BUSCOs (F)|
|63 Missing BUSCOs (M)|
|3285 Total BUSCO groups searched|
=====

```

## 5.2 QUAST

```

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

```

```

1 #Results:
2 cd /home/blowflies/genome_annotation/cmeg/5-RepeatMasker/2-QUAST/cmeg.fa/quast
3 cat report.txt

```



```

All statistics are based on contigs of size >= 3000 bp, unl
Assembly                                cmeg_N_genome_final.fa.masked
# contigs (>= 0 bp)                      760
# contigs (>= 1000 bp)                   760
# contigs (>= 5000 bp)                   760
# contigs (>= 10000 bp)                  760
# contigs (>= 25000 bp)                   703
# contigs (>= 50000 bp)                   601
Total length (>= 0 bp)                   671171543
Total length (>= 1000 bp)                 671171543
Total length (>= 5000 bp)                 671171543
Total length (>= 10000 bp)                671171543
Total length (>= 25000 bp)                670017351
Total length (>= 50000 bp)                666394117
# contigs                                760
Largest contig                           14378473
Total length                             671171543
Estimated reference length                500000000
GC (%)                                   29.14
N50                                       2214294
NG50                                      2991992
N75                                       1082448
NG75                                      2010070
L50                                        84
LG50                                       50
L75                                        190
LG75                                       102
# total reads                            2110
# left                                   0
# right                                  0
Mapped (%)                               100.0
Properly paired (%)                       0.0
Avg. coverage depth                       1
Coverage >= 1x (%)                        99.95
# N's per 100 kbp                         0.05
# predicted rRNA genes                    340 + 125 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/cmeg_pool_RNA /home/blowflies/genome_annotation/cmeg/
3  #Renaming
4  cd /home/blowflies/genome_annotation/cmeg/
5  mv cmeg_pool_RNA 8-cmeg_pool_RNA
6  #checking md5
7  cd 8-cmeg_pool_RNA/
8  cat MD5.txt
9  #9befbbedbbd4e5d9b79a97d29eefc0be  Cmeg_1.fq.gz
10 #1be90c92ec0a0e97240fe0ce21d4a487  Cmeg_2.fq.gz
11 md5sum Cmeg*
12 #9befbbedbbd4e5d9b79a97d29eefc0be  Cmeg_1.fq.gz
13 #1be90c92ec0a0e97240fe0ce21d4a487  Cmeg_2.fq.gz
14 mkdir 0-raw_reads
15 mv *.fq.gz 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://stab.st-andrews.ac.uk/wiki/index.php/FASTQC_and_MultiQC)  
-> <https://multiqc.info/>

```
cd /home/blowflies/genome_annotation/cput/4-cmeg_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
```

Results: [https://drive.google.com/file/d/1ggrh-DgYnWfA8WH2X4tdXRJ4Xo9ZHac/view?usp=share\\_link](https://drive.google.com/file/d/1ggrh-DgYnWfA8WH2X4tdXRJ4Xo9ZHac/view?usp=share_link)

Coping multiqc report to a local computer

```
scp -P 2205 vanessa@143.107.244.181:/home/blowflies/genome_annotation/cmeg/4-cmeg_pool_RNA/1-QC/multiqc_report.html /mnt/c/Users/vansc/Downloads
```

## 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/cmeg/8-cmeg_pool_RNA/0-raw-reads
mkdir ../2-trimming
screen
TrimmomaticPE Cmeg_1.fq.gz Cmeg_2.fq.gz -threads 8 -baseout
/home/blowflies/genome_annotation/cmeg/8-cmeg_pool_RNA/2-trimming/cmeg.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/cmeg/4-cmeg_pool_RNA/2-trimming
# QC
fastqc *.gz

# here -> /home/blowflies/genome_annotation/cmeg/4-cmeg_pool_RNA/1-QC
mkdir trimmed_reads_qc

# in this directory -> /home/blowflies/genome_annotation/cmeg/4-cmeg_pool_RNA/2-trimming
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/cmeg/0-genome
2 |
3 | mv purged.fa cmeg.fa
```

## Quality Control of raw reads

```

1 fastqc /home/cunha/03-RNA/01-Reads/cmeg/sra_ncbi/*fastq
2
3 # in /home/cunha/03-RNA/01-Reads/cmeg/sra_ncbi/02-qc
4 multiqc .

```

## Trimming

```

1 #!/bin/bash
2
3 #SBATCH -N 1-10
4 #SBATCH -n 20
5 #SBATCH -t 240:00:00
6 #SBATCH -p long
7 #SBATCH -o cmeg_trimmomatic.out
8
9 srun /home/cunha/anaconda3/bin/trimmomatic PE SRR1660427_1.fastq SRR1660427_2.fastq
10 srun /home/cunha/anaconda3/bin/trimmomatic PE SRR1663113_1.fastq SRR1663113_2.fastq
11 srun /home/cunha/anaconda3/bin/trimmomatic PE SRR1663114_1.fastq SRR1663114_2.fastq

```

## 7 - Trinity

We had to send the files to /home/cunha/03-RNA/01-Reads/cmeg directory and then rename all of them to make it easier to run trinity

```

1 #sending the files
2 rsync -av /home/cunha/Genotype_Phenotype/0-sequences/1-transcriptome/Cmeg/1-Reads/02-qc/ /home/cunha/03-RNA/01-Reads/cmeg/
3
4 #renaming them
5 # rename files
6 for f in *fastq.gz; do mv -- "$f" "${f%.fastq.gz}.fq.gz"; done

```

## Running Trinity

```

1 #!/bin/bash
2
3 #SBATCH --job-name trinity_cmeg ## nome que aparecerá na fila
4 #SBATCH --output trinity_cmeg.out ## nome do arquivo de saída; o %j é igual a job_id
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 vírgula)
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

```

The final assembly is here: /home/cunha/03-RNA/02-Trinity (Darwin). And it was renamed to cmeg\_trinity.fasta

## Assembly statistics

1	file	format	type	num_seqs	sum_len	min_len
2	cmeg_trinity_all.Trinity.fasta	FASTA	DNA	154,054	137,251,875	166

## BUSCO (transcriptome quality)

```

1 mkdir /home/cunha/03-RNA/03-Busco/BUSCO_RNA_all
2
3 docker run -u $(id -u) -v $(pwd):/busco_wd ezlabgva/busco:v5.4.7_cv1 busco -i

```

## 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/cmeg
7 mkdir 9-STAR_new
8 chmod 777 9-STAR_new
9 cd 9-STAR_new
10 mkdir star_index
11 chmod 777 star_index
12
13 STAR --runThreadN 6 --runMode genomeGenerate --genomeDir /home/blowflies/genome_annotation/cmeg/9-STAR_new/star_index
14
15 #OLD
16 #STAR --runThreadN 6 --runMode genomeGenerate --genomeDir /home/blowflies/genome_annotation/cmeg/9-STAR_new/star_index
```

Before the alignment itself we had to concatenate all the fastq files available. The files are in Rosalind (/home/blowflies/genome\_annotation/cmeg/8-cmeg\_pool\_RNA/2-trimming/all)

```
1 cd /home/blowflies/genome_annotation/cmeg/8-cmeg_pool_RNA/2-trimming/all
2 cat *_1P.fq.gz > cmeg_all_1P.fq.gz
3 cat *_2P.fq.gz > cmeg_all_2P.fq.gz
```

## STAR alignment

```
1 cd /home/blowflies/genome_annotation/cmeg/8-cmeg_pool_RNA/2-trimming/all
2
3 for i in *_1P.fq.gz; do STAR --runMode alignReads --readFilesCommand zcat --
4
5 #OLD
6 #for i in *_1P.fq.gz; do
7 #STAR --runMode alignReads --readFilesCommand zcat --outSAMtype BAM SortedBy
```

Results:

[illegible]

Trying again:

```
1 | #On Rosalind
2 | cd /home/blowflies/genome_annotation/cmeg/8-cmeg_pool_RNA/2-trimming/all
3 | for i in *_1P.fq.gz; do STAR --runMode alignReads --readFilesCommand zcat --
```

**It worked and finished successfully!!!**

## 9- BRAKER3 - Structural annotation

We ran BRAKER3 for structural annotation.

"BRAKER3 is the latest pipeline in the BRAKER suite. It enables the usage of RNA-seq and protein data in a fully automated pipeline to train and predict highly reliable genes with GeneMark-ETP and AUGUSTUS. The result of the pipeline is the combined gene set of both gene prediction tools, which only contains genes with very high support from extrinsic evidence." (<https://github.com/Gaius-Augustus/BRAKER>)

We needed to copied the masked genome and the sorted bam to /home/diniz/programs/braker in Rosalind to run BRAKER3.

```

1 | cd /home/diniz/programs/braker
2 | mkdir Cmeg_braker_new
3 | cd Cmeg_braker_new
4 | cp /home/blowflies/genome_annotation/cmeg/5-RepeatMasker/cmeg_N_genome_final
md5sum cmeg_N_genome_final.fa.masked
5 | #md5 checked
6 |
7 | cp /home/blowflies/genome_annotation/cmeg/9-STAR_new/cmeg_allAligned.out.bam
md5sum cmeg_allAligned.out.bam
8 | #9ffcfde10312b6eca9680911ad6f51ec
9 |
10 | md5sum /home/blowflies/genome_annotation/cmeg/9-STAR_new/cmeg_allAligned.out
11 | #9ffcfde10312b6eca9680911ad6f51ec

```

#OLD

STAR couldn't allocate RAM memory for bam sorting, so we did it manually with samtools

```

1 | cd /home/blowflies/genome_annotation/cmeg/9-STAR_new/
2 | samtools sort -o cmeg_allAligned.sort.out.bam cmeg_allAligned.out.bam

```

After that we moved the genome file Cmeg\_masked.fasta (from /home/pedro/Non\_Coding\_Element\_Evolution/3-Masking/2-RepeatMasker/) and the sorted bam to /home/diniz/braker in rosaling and ran BRAKER3. We ran the command line in a bash script (cmeg\_annot.sh).

## BRAKER3 run

```

1 | # need to run first
2 | export BRAKER_SIF=/home/diniz/programs/braker/braker3.sif
3 | nohup bash cmeg_annot.sh &
4 |
5 | # the content of cmeg_annot.sh
6 | #New
7 | singularity exec /home/diniz/programs/braker/braker3.sif braker.pl --genome=,
8 | #--skipOptimize
9 |
10 |
11 | #OLD
12 | #braker.pl --genome=/home/diniz/braker/cmeg_final_results/Cmeg_masked.fasta

```

## Busco evaluation with protein sequences

```

1 | cd /home/diniz/programs/braker/Cmeg_braker_new/
2 |
3 | docker run -u $(id -u) -v $(pwd):/busco_wd ezlabgva/busco:v5.4.4_cv1 busco -i

```

Results:

```

-----
|Results from dataset diptera_odb10|
-----
|C:97.1%[S:73.4%,D:23.7%],F:0.7%,M:2.2%,n:3285|
|3188 Complete BUSCOs (C)|
|2410 Complete and single-copy BUSCOs (S)|
|778 Complete and duplicated BUSCOs (D)|
|23 Fragmented BUSCOs (F)|
|74 Missing BUSCOs (M)|
|3285 Total BUSCO groups searched|
-----

```

## BRAKER3 second run

```

1 | # need to run first
2 | export BRAKER_SIF=/home/diniz/programs/braker/braker3.sif
3 |
4 | cd /home/diniz/programs/braker/Cmeg_braker_new
5 | mkdir second_run
6 | cd second_run
7 | nano cmeg_annot_2.sh
8 | #cmeg_annot_2.sh
9 | singularity exec /home/diniz/programs/braker/braker3.sif braker.pl --genome=,

```

```

1 | #Run
2 | cd /home/diniz/programs/braker/Cmeg_braker_new/second_run
3 | nohup bash cmeg_annot_2. &

```

How many transcripts

```
1 | cd /home/diniz/programs/braker/Cmeg_braker_new/second_run
2 | grep -c ">" braker.codingseq
3 | #28019
```

**Total transcripts: 28019**

## Busco evaluation with protein sequences from second Braker run

```
1 | cd /home/diniz/programs/braker/Cmeg_braker_new/second_run
2 | mv braker.gtf cmeg_braker.gtf
3 | mv braker.aa cmeg_braker.aa
4 | mv braker.codingseq cmeg_braker.codingseq
5 | md5sum cmeg*
6 | #410ad1b4fc14ebaf5e19f2a1951ca cmeg_braker.aa
7 | #fc8e26a8c1721f6ba954e0e8c5fdbe9f cmeg_braker.codingseq
8 | #5b340dc2a5033183337f48b66aac6926 cmeg_braker.gtf
9 |
10 |
11 | docker run -u $(id -u) -v $(pwd):/busco_wd ezlabgva/busco:v5.4.4_cv1 busco -:
```

**Results:**

```
|Results from dataset diptera_odb10|
|C:97.0%[S:78.6%,D:18.4%],F:0.7%,M:2.3%,n:3285|
|3185 Complete BUSCOs (C)|
|2581 Complete and single-copy BUSCOs (S)|
|604 Complete and duplicated BUSCOs (D)|
|24 Fragmented BUSCOs (F)|
|76 Missing BUSCOs (M)|
|3285 Total BUSCO groups searched|
```

Coping all braker files to /home/blowflies/genome\_annotation/cmeg

```
1 | cd /home/blowflies/genome_annotation/cmeg
2 | mkdir 10-Braker_new
3 | cd 10-Braker_new
4 | sudo cp -r /home/diniz/programs/
5 | braker/Cmeg_braker_new ./
6 | #md5 checked
```

## 11- EnTAP

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

```

1  #We copied the aminoacid file from Rosalind (/home/diniz/programs/braker/Cmeg_braker_new)
2  ssh martins@lem.ib.usp.br -p 4988
3  55rx64yz$
4  cd /home/martins/EnTAP_restart/cmeg/
5  mkdir new
6  cd new
7
8  #copying new file from Rolalind to Darwin server
9  scp -P 2205 diniz@143.107.244.181:/home/diniz/programs/braker/Cmeg_braker_new/cmeg_braker_aa
10 #Checking md5
11 #md5 on Darwin
12 md5sum cmeg_braker.aa
13 #410ad1b4fc14ebalfaf5e19f2a1951ca
14 #md5 on Rosalind
15 #410ad1b4fc14ebalfaf5e19f2a1951ca
16
17 #Running EnTAP
18 #!/bin/bash
19
20 #SBATCH --job-name entap_cmeg ## nome que aparecerá na fila
21 #SBATCH --output entap_cmeg.out ## nome do arquivo de saída; o %j é igual a job-id
22 #SBATCH --ntasks=1 ## número de tarefas (análises) a serem executadas
23 #SBATCH --cpus-per-task=10 ## o número de threads alocados para cada tarefa
24 #SBATCH --mem-per-cpu=1000M # memória por núcleo da CPU
25 #SBATCH --partition=long ## as partições a serem executadas (separadas por vírgula)
26 #SBATCH --error=err
27
28 srun EnTAP --runP -i /home/martins/EnTAP_restart/cmeg/new/cmeg_braker.aa -d /home/martins/EnTAP_restart/cmeg/new/
29
30
31 #checking md5
32 md5sum entap_outfiles/final_results/*
33 #d41d8cd98f00b204e9800998ecf8427e  annotated_contam.faa
34 #9fc1ae9cc0246a8202519b62f903e047  annotated_contam_gene_ontology_terms.tsv
35 #12da35de58c8400536adef688be44687  annotated_contam.tsv
36 #b89d36011f0929530fa3b7b76ce0d1d7  annotated.faa
37 #0e2f85b9015bb174fd47329979650832  annotated_gene_ontology_terms.tsv
38 #4438ac6c503707545294f5b43d14cde7  annotated.tsv
39 #b89d36011f0929530fa3b7b76ce0d1d7  annotated_without_contam.faa
40 #0e2f85b9015bb174fd47329979650832  annotated_without_contam_gene_ontology_terms.tsv
41 #4438ac6c503707545294f5b43d14cde7  annotated_without_contam.tsv
42 #1a8f2d7eac09b46de2802807c0853d2b  entap_results.tsv
43 #1e1bd41573cd19d2704e67bef085d671  unannotated.faa
44 #37ef383804571feb038978980234b231  unannotated.tsv
45
46
47 #Copying the results to Rosalind
48 scp -r -P 2205 entap_outfiles/ vanessa@143.107.244.181:/home/blowflies/genome_annotation/cmeg/11-EnTAP_new/entap_outfiles/final_results/

```

#### Checking md5

```

1  cd /home/blowflies/genome_annotation/cmeg/11-EnTAP_new/entap_outfiles/final_results/
2  md5sum *
3  #d41d8cd98f00b204e9800998ecf8427e  annotated_contam.faa
4  #9fc1ae9cc0246a8202519b62f903e047  annotated_contam_gene_ontology_terms.tsv
5  #12da35de58c8400536adef688be44687  annotated_contam.tsv
6  #b89d36011f0929530fa3b7b76ce0d1d7  annotated.faa
7  #0e2f85b9015bb174fd47329979650832  annotated_gene_ontology_terms.tsv
8  #4438ac6c503707545294f5b43d14cde7  annotated.tsv
9  #b89d36011f0929530fa3b7b76ce0d1d7  annotated_without_contam.faa
10 #0e2f85b9015bb174fd47329979650832  annotated_without_contam_gene_ontology_terms.tsv
11 #4438ac6c503707545294f5b43d14cde7  annotated_without_contam.tsv
12 #1a8f2d7eac09b46de2802807c0853d2b  entap_results.tsv
13 #1e1bd41573cd19d2704e67bef085d671  unannotated.faa
14 #37ef383804571feb038978980234b231  unannotated.tsv
15

```

#### Making an unique gtf file with augustus and ENTAP outputs

files:

- /home/blowflies/genome\_annotation/cmeg/11-EnTAP\_new/entap\_outfiles/final\_results/entap\_results.tsv
- 1a8f2d7eac09b46de2802807c0853d2b
- /home/diniz/programs/braker/Cmeg\_braker\_new/second\_run/cmeg\_braker.gtf

We downloaded the files to a local computer and checked md5.

#md5 checked

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 = "braker.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 = "cmeg_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: 18498 (out of 28019)**

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

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

## 12- Final files

```
1 cd /home/blowflies/genome_annotation/cmeg/12-final_files
2 md5sum *
3 #74d1c8756206130b81c9db811ad7713a Cmeg_annot.gtf
4 #fc8e26a8c1721f6ba954e0e8c5fdbe9f Cmeg_cds.fa
5 #03fcceba9f9f74c8e419382978d184655 Cmeg_genome.fa
6 #410ad1b4fc14ebaf5e19f2a1951ca Cmeg_protein.aa
```

type	original/copy	file	Path	md5
Raw genome	original	XDOVE_20221013_S64018_PL100269092-1_B01.ccs.fastq.gz	/home/Reference_genomes/Cmegacephala	4c8k
Raw genome	original	purged.fa	/home/Reference_genomes/Cmegacephala	3224
Raw genome	copy	cmeg.fa	/home/blowflies/genome_annotation/cmeg/0-genome	3224

type	original/copy	file	Path	md5
Mitochondrial genome	original	cmeg_mit.fa	/home/blowflies/genome_annotation/cmeg/0-genome	faff2f427a88295t
Nuclear genome	original	cmeg_N_genome_final.fa	/home/blowflies/genome_annotation/cmeg/0-genome	657b8451d41d97

type	original/copy	file	Path	md5
RNA-seq	original	Cmeg_1.fq.gz	/home/Raw_seqs/cmeg_pool_RNA	9befbbdbbd4e5d9b7
RNA-seq	original	Cmeg_2.fq.gz	/home/Raw_seqs/cmeg_pool_RNA	1be90c92ec0a0e9724
RNA-seq unzipped unzipped	original	Cmeg_1.fq	/home/Raw_seqs/cmeg_pool_RNA	ed5e4a8587d3e36a4e



RNA-seq unzipped	original	Cmeg_2.fq	/home/Raw_seqs/cmeg_pool_RNA	31dbb84405770c7d79
Trimmed reads zipped	original	cmeg_all_1P.fq.gz	/home/blowflies/genome_annotation/cmeg/8-cmeg_pool_RNA/2-trimming/all	c21045ff1ef48ace62fd
Trimmed reads zipped	original	cmeg_all_2P.fq.gz	/home/blowflies/genome_annotation/cmeg/8-cmeg_pool_RNA/2-trimming/all	8d5c4fe5d822e318dc
Transcriptome	original	cmeg_trinity.fasta	/home/cunha/03-RNA/02-Trinity	ff2f24c065bb88266e7

type	original/copy	file	Path	md5
Masked genome	original	cmeg_N_genome_final.fa.masked	/home/blowflies/genome_annotation/cmeg/5-RepeatMasker	03fccebaf9
Masked genome	old copy	cmeg_masked.fasta	/home/blowflies/genome_annotation/cmeg/9-STAR	a3ead7c04
Masked genome	old copy	Cmeg_masked.fasta	/home/blowflies/genome_annotation/cmeg/10-BRAKER3	a3ead7c04
Masked genome	old copy	Cmeg_masked.fasta	/home/diniz/programs/braker/Cmeg_2/10-BRAKER3	a3ead7c04
Masked genome	copy	Cmeg_masked.fasta	/home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker	03fccebaf9
Masked genome	new copy	cmeg_N_genome_final.fa.masked	/home/diniz/programs/braker/Cmeg_braker_new	03fccebaf9

type	original/copy	file	Path	md5
Proteome	original	cmeg_braker.aa	/home/diniz/programs/braker/Cmeg_braker_new/second_run	410ad1b4fc14
Proteome	copy	cmeg_braker.aa	/home/martins/EnTAP_restart/Proteomes	
Proteome	copy	cmeg_braker.aa	/home/00-Sequences/Cmegacephala/01-Genomic_data/2023	20e2e610ee6
Proteome	copy	cmeg_braker.aa	/home/blowflies/genome_annotation/cmeg/10-Braker_new/Cmeg_braker_new/second_run	410ad1b4fc14

type	original/copy	file	Path	md5
gtf output Braker3 second run	original	cmeg_braker.gtf	/home/diniz/programs/braker/Cmeg_braker_new/second_run	5b340dc2a503
gtf output Braker3 second run	copy	cmeg_braker.gtf	/home/blowflies/genome_annotation/cmeg/10-Braker_new/Cmeg_braker_new/second_run	5b340dc2a503
gtf output Braker3 second run	copy		Computador do Diniz	
Final gtf	original		Computador do Diniz	
Final gtf	copy		Rosalind	

type	original/copy	file	Path	md5

Condensingseq	original	braker.codingseq	/home/diniz/programs/braker/Cmeg_braker_new/second_run	fc8e
Condensingseq	copy	cmeg_braker.codingseq	/home/blowflies/genome_annotation/cmeg/10-Braker_new/Cmeg_braker_new/second_run	fc8e
Condensingseq	copy	cmeg_braker.codingseq	/home/00-Sequences/Cmegacephala/01-Genomic_data/2023	f85:

**FALTA: conferir md5 dos arquivos usados no script do R do entap e cmeg\_entap\_final.gtf que está no computador do Diniz**

#### PROBLEMAS:

- arquivo output do RepeatMasker diferente dos demais que usamos de input nos outros programas. Rodar tudo novamente??
- arquivo proteoma da pasta 00-sequences está com md5 diferente do original
- arquivo codingseq da pasta 00-sequences está com md5 diferente do original
- arquivo gtf que está na pasta 00-sequences não é o final

## OBSOLETE

## 10- gFACs

<https://gfacs.readthedocs.io/en/latest/Flags/index.html>

```
1 | cd /home/blowflies/genome_annotation/cmeg/
2 | mkdir 11-gfacs
3 | cd /home/blowflies/genome_annotation/cmeg/11-gfacs
4 | mkdir results
5 | sudo perl /gFACs-master/gFACs.pl -f braker_2.1.2_gtf -p cmeg --rem-all-incom
```

Results:

Number of genes (Augustus/BRAKER): 29938

Number of genes (gFACs): 29258

## 11- EnTAP

We did it on Darwin

```
1 | cd /home/martins/EnTAP
2 |
3 | # The gFACs outputs for all species are here (*_genes.fasta.faa)
4 | mkdir Proteomes
5 |
6 | mkdir cmeg
7 | cd cmeg
8 | EnTAP --runP -i /home/martins/EnTAP/Proteomes/cmeg_genes.fasta.faa -d /home/r
```

## 12. gFACs again (with EnTAP output)

```
1 | cd /home/blowflies/genome_annotation/cmeg/
2 | mkdir 13-gfacs_entap
3 | cd /home/blowflies/genome_annotation/cmeg/13-gfacs_entap
4 | mkdir results
5 | sudo perl /gFACs-master/gFACs.pl -f gFACs_gene_table -p cmeg --rem-all-incom
```

## 13. Final annotation

Final files are here:

```
1 | /dados/home/blowflies/genome_annotation/cmeg/14-final_annot
2 | # 24543 gene models
```

Final busco

```
1 | docker run -u $(id -u) -v $(pwd):/busco_wd ezlabgva/busco:v5.4.4_cv1 busco -:
```

```
***** Results: *****  
  
C:96.5%[S:79.3%,D:17.2%],F:0.9%,M:2.6%,n:3285  
3169   Complete BUSCOs (C)  
2605   Complete and single-copy BUSCOs (S)  
564    Complete and duplicated BUSCOs (D)  
29     Fragmented BUSCOs (F)  
87     Missing BUSCOs (M)  
3285   Total BUSCO groups searched
```

## OBSOLETE STUFF

### RepeatModeler

```
1 | mkdir 4-RepeatModeler  
2 | cp /home/Reference_genomes/Cmegacephala/purged.fa ./cmeg.fa  
3 |  
4 | #Database  
5 | /RepeatModeler-2.0.4/BuildDatabase -name cmeg_database cmeg.fa  
6 |  
7 | #RepeatModeler  
8 | screen  
9 | /RepeatModeler-2.0.4/RepeatModeler -database cmeg_database -threads 20 -LTRS  
10 | #[1] 46887
```

### RNA reads from NCBI

We downloaded RNA SRAs available from NCBI to use with our own RNA-seq to use it as more evidence for our structural annotation.

We downloaded the SRAs with SRA Toolkit in Darwin. We used a .txt list with the accession numbers.

cmeg\_list.txt

SRR1660427

SRR1663113

SRR1663114

They were downloaded in /home/cunha/03-RNA/01-Reads/cmeg/sra\_ncbi:

```
1 | prefetch --option-file cmeg_list.txt
```

And then we had to convert the sequences to fastq:

```
1 | fasterq-dump --split-files SRR1660427.sra  
2 | fasterq-dump --split-files SRR1663113.sra  
3 | fasterq-dump --split-files SRR1663114.sra
```

The SRA files are in /home/cunha/03-RNA/01-Reads/cmeg/sra\_ncbi/01-sra\_files

## Trinity (first)

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/cmeg/8-cmeg_pool_RNA/2-trimming
```

### Assembling the transcriptome

```

1 | cd /home/cunha/03-RNA/01-Reads/cmeg
2 |
3 | docker run --rm -v`pwd`:`pwd` trinityrnaseq/trinityrnaseq Trinity --seqType :
4 |
5 | mv cmeg_trinity.* /home/cunha/03-RNA/02-Trinity/cmeg_trinity
6 |
7 | # Transcriptome size
8 | grep -c ">" cmeg_trinity.Trinity.fasta # 54372

```

## 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:v5.4.4

```

Complete (all)	Complete Single	Complete Dup.	Fragmented	Missing
74.8	36.2	38.6	6.9	18.3

## 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/](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 megacephala]/' Cmeg_genome.fa
2 | singularity run agat_1.0.0--pl5321hd78af_0.sif
3 | agat_convert_sp_gxf2gxf.pl --gtf Cmeg_annot.gtf --output Cmeg_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 Cmeg_annot.gff \
7 | -i Cmeg_genome.fa \
8 | -locus-tag-prefix Cmeg \
9 | -o cmeg.sqn \
10 | -Z \
11 | -V b

```