

Log: Genome Annotation *Lucilia eximia*

tags: `Genome Annotation` `Lexi`

Table of Contents

- [Log: Genome Annotation *Lucilia eximia*](#)
 - [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](#)
- [4- RepeatModeler](#)
 - [Installing with docker](#)
 - [Creating database](#)
 - [Running](#)
 - [Copying the results](#)
- [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](#)
- [7- Trinity](#)
 - [Running Trinity](#)
 - [Assembly statistics](#)
 - [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](#)
- [OBSOLETE](#)
- [10- gFACs](#)
- [11- EnTAP](#)
- [12. gFACs again \(with EnTAP output\)](#)
- [13. Final annotation](#)
- [OBSOLETE STUFF](#)
 - [Trinity \(first\)](#)

- Moving trimmed reads to Darwin
- Assembling the transcriptome
- BUSCO (transcriptome quality)
- 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 2723669

Coverage (x) 65

HPA Length (bp) 1190603997

HPA N50 (bp) 2920813

HPA N90 (bp) 507262

HPA L50 118

HPA L90 467

FA Length (bp) 579241898

FA N50 (bp) 4214109

FA N90 (bp) 1077085

FA L50 42

FA L90 140

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/lexi
3
4 # new directory
5 mkdir 0-genome
6 cd 0-genome
7
8 cp /home/Reference_genomes/Leximia/purged.fa .
9 cp /home/Reference_genomes/Leximia/XDOVE_20221110_S64411e_PL100270436-1_C01.c
```

1- BUSCO First

Busco version 5.3.2

```
1 mkdir /home/blowflies/genome_annotation/lexi/1-busco_first
2 cd /home/blowflies/genome_annotation/lexi/1-busco_first
3
4 #sudo docker pull ezlabgva/busco:v5.4.4_cv1
5
6 cp ../0-genome/lexi.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:98.9%[S:65.2%,D:33.7%],F:0.3%,M:0.8%,n:3285|
|3250   Complete BUSCOs (C)                  |
|2142   Complete and single-copy BUSCOs (S)   |
|1108   Complete and duplicated BUSCOs (D)     |
|9       Fragmented BUSCOs (F)                |
|26      Missing BUSCOs (M)                   |
|3285   Total BUSCO groups searched           |
|-----|
```

2- QUAST

```
1 # we are here
2 /home/blowflies/genome_annotation/lexi
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

We had to use the NCBI's reference mitochondrial genome of *Lucilia cuprina* since there's not a reference mitochondrial genome for *Lucilia eximia*.

https://www.ncbi.nlm.nih.gov/nuccore/NC_002660.1

search terms: "lucilia cuprina [ORGN] AND mitochondrial genome"

NC_019573.1

```
1 | esearch -db nuccore -query NC_019573.1 | efetch -format fasta > lcup_mitocho
2 | cd home/blowflies/genome_annotation/lexi/3-mitochondrial_genome
3 | mkdir 1-blastn
```

Blastn

version 2.11

```
1 | cd home/blowflies/genome_annotation/lexi/3-mitochondrial_genome/1-blastn
2 |
3 | #Making database using the genome
4 | makeblastdb -in ../lcup_mitochondrial_genome.fa -dbtype nucl -out lexi_mit_db
5 |
6 | #Running
7 | blastn -task blastn -evaluate 0.00001 -db ../lexi_mit_database -query ../../0-genome.fa
8 |
9 | #extract mitochondrial genome (scaffold with the bigger alignement with the mitochondria)
10 | seqkit grep -p "ptg0010921" lexi.fa > ptg0010921.fa
11 |
12 | #renaming
13 | mv ptg0010921.fa lexi_mit_scaff.fa
14 |
15 | #removing the mitochondria from the genome using grep invert matching
16 | seqkit grep -v -p "ptg0010921" lexi.fa > lexi_N_genome.fa
17 |
18 | #checking the sequences
19 | grep -c ">" lexi*
20 | #lexi.fa:369
21 | #lexi_mit_scaff.fa:1
22 | #lexi_N_genome.fa:368
```

Then we used the blast output table to

```
1 | # removing duplicates - getting a new fasta from positions 13326:29418
2 | seqkit subseq lexi_mit_scaff.fa -r 13326:29418 > lexi_mit_nodup.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/lexi/3-mitochondrial_genome/2-MITOS_results).

The job settings were:

Job ID: lexi_mit_genome

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

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: fd973780a3fd
5 | container-name: friendly_archimedes
6 | container-image: dfam/tetools:latest
7 |
8 | #To attach the container and continue running press CTRL+P, then CTRL+Q
```

```
(base) cunha@darwin:~/00-Genomes$ docker ps
CONTAINER ID   IMAGE          COMMAND                  CREATED        STATUS        PORTS        NAMES
fd973780a3fd   dfam/tetools:latest   "bash"                  22 minutes ago   Up 14 minutes           friendly_archimedes
```

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 ./lexi_N_genome.fa friendly_archimedes:/home
4 |
5 | # Getting inside the container
6 | docker exec -it friendly_archimedes /bin/bash/
7 |
8 | #Database
9 | BuildDatabase -name lexi_database lexi_N_genome.fa
```

Running

```
1 | RepeatModeler -database lexi_database -threads 20 -LTRStruct >log 2>err
```

Copying the results

```
1 | cd /home/cunha/01-RepeatModeler/lexi
2 |
3 | docker cp -a friendly_archimedes:/home/lexi/lexi_database-families.stk ./
4 | docker cp -a friendly_archimedes:/home/lexi/lexi_database-families.fa ./
5 | docker cp -a friendly_archimedes:/home/lexi/lexi_database-rmod.log ./
```

5- RepeatMasker

Installing steps can be found in the *Chrysomya megacephala* log

(https://hackmd.io/ziOztK1MQZecVm0_qDAr1g)

```

1 | # we start here
2 | /home/cunha/02-RepeatMasker/lexi
3 |
4 | # running
5 | RepeatMasker -lib /home/cunha/01-RepeatModeler/lexi/lexi_database-families.f

```

2.5.1 BUSCO

```

1 | cd /home/blowflies/genome_annotation/lexi/5-RepeatMasker
2 | mkdir 1-BUSCO
3 | cp /home/blowflies/genome_annotation/lexi/9-STAR/Lexi_masked.fasta 1-BUSCO/
4 | cd 1-BUSCO
5 | md5sum Lexi_masked.fasta
6 | sudo docker run -u $(id -u) -v $(pwd):/busco_wd ezlabgva/busco:v5.4.7_cv1 bu

```

#Results

```

-----
|Results from dataset diptera_odb10|
-----
|C:99.0%[S:65.5%,D:33.5%],F:0.3%,M:0.7%,n:3285|
|3251 Complete BUSCOs (C)|
|2151 Complete and single-copy BUSCOs (S)|
|1100 Complete and duplicated BUSCOs (D)|
|10 Fragmented BUSCOs (F)|
|24 Missing BUSCOs (M)|
|3285 Total BUSCO groups searched|
-----

```

2.5.2 QUAST

```

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

```

```

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

```

```

Assembly Lexi_masked
# contigs (>= 0 bp) 368
# contigs (>= 1000 bp) 368
# contigs (>= 5000 bp) 368
# contigs (>= 10000 bp) 368
# contigs (>= 25000 bp) 325
# contigs (>= 50000 bp) 306
Total length (>= 0 bp) 579203135
Total length (>= 1000 bp) 579203135
Total length (>= 5000 bp) 579203135
Total length (>= 10000 bp) 579203135
Total length (>= 25000 bp) 578414946
Total length (>= 50000 bp) 577733876
# contigs 368
Largest contig 15451354
Total length 579203135
Estimated reference length 500000000
GC (%) 29.30
N50 4214109
NG50 4760762
N75 2285326
NG75 3010171
L50 42
LG50 33
L75 88
LG75 66
# total reads 729
# left 0
# right 0
Mapped (%) 100.0
Properly paired (%) 0.0
Avg. coverage depth 1
Coverage >= 1x (%) 99.99
# N's per 100 kbp 0.06
# predicted rRNA genes 182 + 17 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/lexi_pool_RNA /home/blowflies/genome_annotation/lexi/
3 #Renaming
4 cd /home/blowflies/genome_annotation/lexi/
5 mv lexi_pool_RNA 8-lexi_pool_RNA
6 #checking md5
7 cd 8-lexi_pool_RNA/
8 cat MD5.txt
9 #856d439819ae45beefd5e4c5785bd4a0 Lexi_1.fq.gz
10 #cddb764449f3096589f36afbc0bad1a3 Lexi_2.fq.gz
11 md5sum Lexi*
12 #856d439819ae45beefd5e4c5785bd4a0 Lexi_1.fq.gz
13 #cddb764449f3096589f36afbc0bad1a3 Lexi_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://multiqc.info/>

```
cd /home/blowflies/genome_annotation/lexi/4-lexi_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/lexi/4-lexi_pool_RNA/1-QC/multiqc_report.html /mnt/c/Users/vansc/Downloads
```

Results:

https://drive.google.com/file/d/11K0LqILvPrOrGrpvFY7DZRB5OosV_vzs/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/lexi/4-lexi_pool_RNA/0-raw-reads
mkdir ../2-trimming
screen
TrimmomaticPE Lexi_1.fq.gz Lexi_2.fq.gz -threads 8 -baseout
/home/blowflies/genome_annotation/lexi/8-lexi_pool_RNA/2-
trimming/lexi.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/lexi/4-lexi_pool_RNA/2-trimming
# QC
fastqc *.gz

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

7- Trinity

We had to send the files from Rosalind to Darwin then rename all of them to make it easier to run trinity

```
1 | #renaming them
2 | # rename files
3 | 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_lexi ## nome que aparecerá na fila
4 | #SBATCH --output trinity_lexi.out ## nome do arquivo de saída; o %j é igual ao 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 lexi_trinity.fasta (md5checked)

Assembly statistics

file	format	type	num_seqs	sum_len	min_len	avg_len
trinity_all.Trinity.fasta	FASTA	DNA	136,730	139,151,237	176	1,032

BUSCO (transcriptome quality)


```

1 | mkdir /home/cunha/03-RNA/03-Busco/BUSCO_RNA_all
2 |
3 | docker pull ezlabgva/busco:v5.4.7_cv1 # just because we didn't have busco on
4 |
5 | docker run -u $(id -u) -v $(pwd):/busco_wd ezlabgva/busco:v5.4.7_cv1 busco -:

```

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/lexi
7 | mkdir 8-STAR
8 | chmod 777 8-STAR
9 | cd 8-STAR
10 | mkdir star_index
11 | chmod 777 star_index
12 |
13 | STAR --runThreadN 8 --runMode genomeGenerate --genomeDir /home/blowflies/geno

```

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

```

1 | cat *_1P.fq.gz > lexi_all_1P.fq.gz
2 | cat *_2P.fq.gz > lexi_all_2P.fq.gz

```

STAR alignent

```

1 | cd /home/blowflies/genome_annotation/lexi/8-lexi_pool_RNA/2-trimming/all
2 |
3 | for i in *_1P.fq.gz; do
4 | STAR --runMode alignReads --readFilesCommand zcat --outSAMtype BAM SortedByCo
5 |
6 | #started at: Feb 10 18:48:37
7 | #finished at: Feb 10 19:17:44

```

9- Braker3

All information is in the Calb file (https://hackmd.io/KMzFBC2aQ9qliTy11f_QhQ)

Final outputs are here: /home/blowflies/genome_annotation/lexi/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 Lexi_2
7 | cp /home/blowflies/genome_annotation/lexi/10-BRAKER3 /home/diniz/programs/br
8 | cd /home/diniz/programs/braker/Lexi_2/10-BRAKER3
9 | mkdir restart
10 | # Run
11 | singularity exec /home/diniz/programs/braker/braker3.sif braker.pl --genome=,

```

Total transcripts: 28884

Busco evaluation with protein sequences from restart Braker run

```

1 | cd /home/diniz/programs/braker/Lexi_2/10-BRAKER3/restart
2 | mv braker.gtf lexi_braker.gtf
3 | mv braker.aa lexi_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: 3285)
# Summarized benchmarking in BUSCO notation for file /busco_wd/lexi_braker.aa
# BUSCO was run in mode: proteins

***** Results: *****

C:98.0%[S:56.5%,D:41.5%],F:0.9%,M:1.1%,n:3285
3217   Complete BUSCOs (C)
1855   Complete and single-copy BUSCOs (S)
1362   Complete and duplicated BUSCOs (D)
29     Fragmented BUSCOs (F)
39     Missing BUSCOs (M)
3285   Total BUSCO groups searched

```

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/Lexi_2/10-BRAKER3/restart/lexi_braker.aa)
2 | cd /home/martins/EnTAP_restart
3 |
4 | EnTAP --runP -i /home/martins/EnTAP_restart/Proteomes/lexi_braker.aa -d /home/martins/EnTAP_restart/Proteomes/
5 |
6 | mv entap_outfiles/ lexi/
7 |
8 | scp -r -P 2205 entap_outfiles/ pedro@143.107.244.181:/home/blowflies/genome_annotation/lexi/11-EnTAP/

```

Checking md5

```

1 | cd /home/blowflies/genome_annotation/lexi/11-EnTAP/entap_outfiles/final_results
2 | md5sum entap_results.tsv
3 | #b520b5c11281d98e4c3919f68e1316e4 entap_results.tsv
4 | cd /home/diniz/programs/braker/Lexi_2/10-BRAKER3/restart
5 | md5sum lexi_braker.gtf
6 | #08af39a581366b0c871fcaa3bbff7e7c lexi_braker.gtf

```

Making an unique gtf file with augustus and ENTAP outputs

files:

- /home/blowflies/genome_annotation/lexi/11-EnTAP/entap_outfiles/final_results/entap_results.tsv
#b520b5c11281d98e4c3919f68e1316e4 entap_results.tsv
- /home/diniz/programs/braker/Lexi_2/10-BRAKER3/restart/lexi_braker.gtf
#08af39a581366b0c871fcaa3bbff7e7c

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 = "lexi_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 = "lexi_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: 21835 (out of 28884)

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

```
1 | cd /Users/diniz/Desktop
2 | md5 Lexi_annot.gtf
3 | #aef125d4a204649fc62f08a36c66716a
4 | scp -P 2205 Lexi_annot.gtf diniz@143.107.244.181:/home/blowflies/genome_anno
5 |
6 | #On Rosalind:
7 | cd /home/blowflies/genome_annotation/lexi/12-final_files/
8 | md5sum Lexi_annot.gtf
9 | #aef125d4a204649fc62f08a36c66716a
```

12-Final files

```
1 | cd /home/blowflies/genome_annotation/lexi/12-final_files
2 | #aef125d4a204649fc62f08a36c66716a Lexi_annot.gtf
3 | #6e24cf78ab05fd2947eecb89f6eb85b6 Lexi_cds.fa
4 | #4117a959ee8e7f397e3648cf53e8153e Lexi_genome.fa
5 | #a4df9e2040357393e1f62715a0ccf4d1 Lexi_protein.aa
6 |
```

type	original/copy	file	Path	md5
Raw genome	original	XDOVE_20221110_S64411e_PL100270436-1_C01.ccs.fastq.gz	/home/Reference_genomes/Leximia	5a7e8
Raw genome	original	purged.fa	/home/Reference_genomes/Leximia	5f5dc
Raw genome	copy	lexi.fa	/home/blowflies/genome_annotation/lexi/0-genome	5f5dc

type	original/copy	file	Path	md5
Mitochondrial genome	original	lexi_mit_nodup.fa	/home/blowflies/genome_annotation/lexi/0-genome	faff2f427a88295b14b4bdf
Nuclear genome	original	lexi_N_genome.fa	/home/blowflies/genome_annotation/lexi/0-genome	b7241aa3c751fc7a1c66df1

type	original/copy	file	Path	md5
RNA-seq	original	Lexi_1.fq.gz	/home/Raw_seqs/lexi_pool_RNA	856d439819ae45beefd5e
RNA-seq	original	Lexi_2.fq.gz	/home/Raw_seqs/lexi_pool_RNA	cdbb764449f3096589f36
RNA-seq unzipped unzipped	original	Lexi_1.fq	/home/Raw_seqs/lexi_pool_RNA	36cb7c42b1e8784c4678b
RNA-seq unzipped	original	Lexi_2.fq	/home/Raw_seqs/lexi_pool_RNA	b68d49d023625ba791b01
Trimmed reads zipped	original	lexi_all_1P.fq.gz	/home/blowflies/genome_annotation/lexi/8-lexi_pool_RNA/2-trimming/all	a3013b0d773c2f49dbc43
Trimmed reads zipped	original	lexi_all_2P.fq.gz	/home/blowflies/genome_annotation/lexi/8-lexi_pool_RNA/2-trimming/all	958d073075ad63426e3da
Transcriptome	original	lexi_trinity.fasta	/home/cunha/03-RNA/02-Trinity	404eecbe03bd747996b48

type	original/copy	file	Path	md5
Masked genome	original	lexi_N_genome.fa.masked	/home/cunha/02-RepeatMasker/lexi	4117a959ee8e7f39
Masked genome	copy	Lexi_masked.fasta	/home/blowflies/genome_annotation/lexi/9-STAR	4117a959ee8e7f39

Masked genome	copy	Lexi_masked.fasta	/home/blowflies/genome_annotation/lexi/10-BRAKER3	4117a959ee8e7f39
Masked genome	copy	Lexi_masked.fasta	/home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker	4117a959ee8e7f39

type	original/copy	file	Path	md5
Proteome	original	lexi_braker.aa	/home/diniz/programs/braker/Lexi_2/10-BRAKER3/restart	a4df9e2040357393e1f62715a0ccf4
Proteome	copy	lexi_braker.aa	/home/martins/EnTAP_restart/Proteomes	a4df9e2040357393e1f62715a0ccf4
Proteome	copy	lexi_braker.aa	/home/00-Sequences/Leximia/01-Genomic_data/2023	bd5a03ebc9b56a26e9fb15921dc514

type	original/copy	file	Path	md5
gtf output Braker3 second run	original	lexi_braker.gtf	/home/diniz/programs/braker/Lexi_2/10-BRAKER3/restart	08af39a581366b0c871fcaa3bbff7e7c
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
Condingseq	original	braker.codingseq	/home/diniz/programs/braker/Lexi_2/10-BRAKER3/restart	6e24cf78ab05fd2947eecb
Condingseq	copy	lexi_braker.codingseq	/home/00-Sequences/Leximia/01-Genomic_data/2023	6cda82d7deb70bf4b6bf17

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 gtf que está na pasta 00-sequences não é o final
- arquivo codingseq da pasta 00-sequences está com md5 diferente do original
- arquivo proteoma da pasta 00-sequences está com md5 diferente do original

OBSOLETE

10- gFACs

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

```

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

Results:
Number of genes (Augustus/BRAKER): 29174

Number of genes (gFACs): 28339

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 lexi
7 | cd lexi
8 | EnTAP --runP -i /home/martins/EnTAP/Proteomes/lexi_genes.fasta.faa -d /home/r
```

I copied the output folder into Rosalind

```
1 | cd /home/blowflies/genome_annotation/lexi
2 |
3 | mkdir 12-EnTAP
4 | cd 12-EnTAP
5 | scp -r -P 4988 martins@lem.ib.usp.br:/home/martins/EnTAP/lexi/entap_outfiles,
```

12. gFACs again (with EnTAP output)

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

13. Final annotation

Final files are here:

```
1 | /dados/home/blowflies/genome_annotation/lexi/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:97.6%[S:57.6%,D:40.0%],F:0.8%,M:1.6%,n:3285
3206 Complete BUSCOs (C)
1891 Complete and single-copy BUSCOs (S)
1315 Complete and duplicated BUSCOs (D)
25 Fragmented BUSCOs (F)
54 Missing BUSCOs (M)
3285 Total BUSCO groups searched
```

OBSOLETE STUFF

Trinity (first)

Transcriptome assembly

Moving trimmed reads to Darwin

```
1 | mkdir /home/cunha/03-RNA/01-Reads/lexi
2 |
3 | scp -P 4988 /home/blowflies/genome_annotation/lexi/8-lexi_pool_RNA/2-trimmin
```

Assembling the transcriptome

```

1 | cd /home/cunha/03-RNA/01-Reads/lexi
2 |
3 | docker run --rm -v `pwd`:`pwd` trinityrnaseq/trinityrnaseq Trinity --seqType
4 |
5 | mv lexi_trinity.* /home/cunha/03-RNA/02-Trinity/lexi_trinity
6 |
7 | # Transcriptome size
8 | grep -c ">" lexi_trinity.Trinity.fasta # 36184

```

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

```

Complete (all)	Complete Single	Complete Dup.	Fragmented	Missing
80	37.4	42.6	3.4	16.6

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=Lucilia eximia]/' Lexi_genome.fa
2 | singularity run agat_1.0.0--pl5321hdfd78af_0.sif
3 | agat_convert_sp_gxf2gxf.pl --gtf Lexi_annot.gtf --output Lexi_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 Lexi_annot.gff \
7 | -i Lexi_genome.fa \
8 | -locus-tag-prefix Lexi \
9 | -o lexi.sqn \
10 | -Z \
11 | -V b

```