

Blowfly non-coding elements

This notebook is for the description of the steps taken for the comparison of non-coding elements in the genomes of Calliphoridae species with different feeding habits.

Table of Contents

- [Blowfly non-coding elements](#)
 - [Table of Contents](#)
 - [Programs used \(and prerequisites\):](#)
 - [0. New Genomes](#)
 - [1. Previously available Genomes](#)
 - [2. Quality of genomes](#)
 - [2.1. BUSCO](#)
 - [2.2. QUAST](#)
 - [Decisions](#)
 - [3. Repeats and soft masking](#)
 - [3.1. RepeatModeler](#)
 - [3.2. RepeatMasker](#)
 - [LTR - Lser](#)
 - [3.3. Summarizing information](#)
 - [3.4. Compressing unmasked genomes](#)
 - [Note on the new genomes](#)
 - [4. Whole genome alignment \(WGA\)](#)
 - [4.1. Installation](#)
 - [4.2. Tree](#)
 - [4.3. Running Cactus](#)
 - [5. PhastCons](#)
 - [5.1. Installation](#)
 - [5.2. Running](#)
 - [5.3. Element filtering](#)
 - [6. Element processing](#)
 - [6.1. Sequences in the other species](#)
 - [6.2. Filtering fasta files](#)
 - [6.3. Alignments](#)
 - [7. PhyloACC](#)
 - [7.1. Installing](#)
 - [7.2. Running](#)
 - [8. Closest genes](#)
 - [8.1. Bed files](#)
 - [8.2. Transcription start sites](#)
 - [8.3. Closest genes](#)
 - [9. Transcription factors](#)
 - [9.1. Getting the CNEEs](#)
 - [9.2. TF binding motifs](#)
 - [9.3. GO enrichment for genes with the TF binding motifs](#)

Programs used (and prerequisites):

Software	Version
----------	---------

BUSCO	5.4.4
QUAST	5.0.2
RepeatMasker	4.1.5
RMblast	2.11.0
TRF	4.09.1
Cactus	2.6.4
HDF5	1.10.1
HAL	2.2
Clapack	3.2.1
PHAST	1.6
bedtools	2.30.0
BEDOPS	2.4.41
Mafft	7.310
PhyloACC	2.2.0
ASTRAL	5.7.1
IQTree	2.0.6
MEME Suite	5.5.4

0. New Genomes

Three new genomes:

Species	Company	Technology	HackMD link
<i>Lucilia eximia</i>	Dovetail	PacBio	link
<i>Chrysomya putoria</i>	Dovetail	PacBio	link
<i>Chrysomya megacephala</i>	Dovetail	PacBio	link

1. Previously available Genomes

Genomes of blowflies were downloaded from the NCBI database. If a RefSeq version were available, it was preferred over the GenBank version. If not, I used the GenBank version.

Species	Accession code	Ref
<i>Lucilia cuprina</i>	GCF_022045245.1	Anstead et al. 2015
<i>Lucilia sericata</i>	GCF_015586225.1	Davis et al. 2021
<i>Chrysomya rufifacies</i> ¹	GCA_014858695.1	Andere et al. 2020
<i>Phormia regina</i>	GCA_001735545.1	Andere et al. 2016
<i>Calliphora vicina</i>	GCA_001017275.1	Vicoso & Bachtrog, 2015
<i>Calliphora vomitoria</i>	GCA_942486065.1	Darwin Tree of Life Project (https://portal.darwintreeoflife.org/data/vomitoria)
<i>Protocalliphora azurea</i>	GCA_932274085.1	Darwin Tree of Life Project (https://portal.darwintreeoflife.org/data/azurea)
<i>Bellardia pandia</i>	GCA_916048285.2	Darwin Tree of Life Project (https://portal.darwintreeoflife.org/data/pandia)

<i>Cochliomyia hominivorax</i> ²	https://datadryad.org/stash/dataset/doi:10.5061/dryad.d7wm37q4j	Tandonnet et al. 2023
<i>Calliphora (Aldrichina) grahami</i> ³	http://dx.doi.org/10.5524/100673	Meng et al. 2020

¹There are three assemblies (arrhenogenic female, thelygenic female, and male). This is the male assembly, I chose it for the Y chromosome.

²There is a GenBank assembly for this species (GCA_004302925.1), but I used the updated version of the genome (Sophie's paper).

³ Not on NCBI, but another repository.

Darwin Tree of Life Project Genomes: <https://portal.darwintreeoflife.org/tracking>
<https://projects.ensembl.org/darwin-tree-of-life/>
<https://portal.darwintreeoflife.org/tree>

I edited sequence ids (left just the first identifier, erasing everything after the first space) to make it shorter and easier to read. These were used afterwards. All original genomes are stored in a compressed file.

```
1 | # Original
2 | /home/pedro/Non_Coding_Element_Evolution/0-Original_Genomes.tar.gz
3 | # Edited
4 | cd /home/pedro/Non_Coding_Element_Evolution/1-Genomes
5 | for i in *; do sed -i 's/ .*//g' $i; done
```

Number of scaffolds/contigs in each assembly:

- Agra.fasta** - 1604
- Bpan.fasta** - 71
- Chom.fasta** - 522
- Cruf.fasta** - 109329
- Cvic.fasta** - 197510
- Cvom.fasta** - 129
- Lcup.fasta** - 8457
- Lser.fasta** - 4371
- Pazu.fasta** - 24
- Preg.fasta** - 192460

Edited genomes were also compressed after the RepeatMasker run finished (see section 3.2.)

2. Quality of genomes

```
1 | mkdir /home/pedro/Non_Coding_Element_Evolution/2-Genome_quality
2 |
3 | cd /home/pedro/Non_Coding_Element_Evolution/2-Genome_quality
```

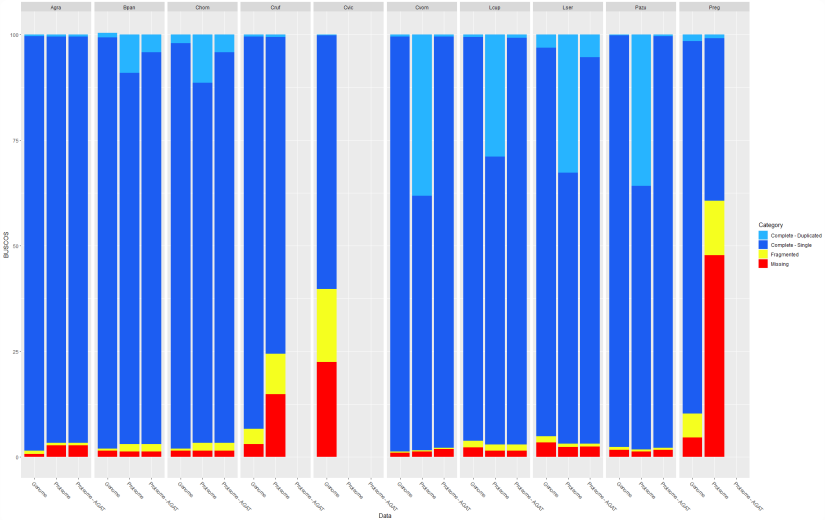
2.1. BUSCO

I ran BUSCO to check for the completeness of each genome.

```
1 | # BUSCO v5.4.4
2 | mkdir /home/pedro/Non_Coding_Element_Evolution/2-Genome_quality/1-BUSCO
3 |
4 | cd 1-BUSCO/
5 |
6 | # For some reason, it only works if the genomes are in the same folder, so I
7 | cp /home/pedro/Non_Coding_Element_Evolution/1-Genomes/*.fasta .
8 |
9 | sudo docker pull ezlabgva/busco:v5.4.4_cv1
10 |
11 | for GENOME in *fasta; do sudo docker run -u $(id -u) -v $(pwd):/busco_wd ezlabgva/busco:v5.4.4_cv1
```

Species	Complete (S + D)	Fragmented	Missing
Agra	98.6 (98.2 + 0.4)	0.7	0.7
Bpan	98.5 (97.4 + 1.1)	0.3	1.2

Chom	98.1 (96.0 + 2.1)	0.5	1.4
Cruf	93.4 (92.9 + 0.5)	3.6	3.0
Cvic	60.3 (60.1 + 0.2)	17.3	22.4
Cvom	98.7 (98.2 + 0.5)	0.3	1.0
Lcup	96.2 (95.6 + 0.6)	1.6	2.2
Lser	95.1 (92.0 + 3.1)	1.5	3.4
Pazu	97.7 (97.5 + 0.2)	0.2	1.6
Preg	89.8 (88.2 + 1.6)	5.6	4.6



OBS: This has Carol's results with the proteomes too, for comparison

2.2. QUAST

```
1 # QUAST v5.0.2
2 mkdir /home/pedro/Non_Coding_Element_Evolution/2-Genome_quality/2-QUAST
3
4 cd 2-QUAST/
5
6 # Running
7 for GENOME in /home/pedro/Non_Coding_Element_Evolution/2-Genome_quality/1-BUSCO:
8
9
10
11 for i in $(ls); do sudo quast.py $i -t 10 --eukaryote --large --memory-efficient
```

Species	#contigs	#contigs>50kb	Length	Length>50kb	N ₅₀	N ₇₅	L ₅₀	L ₇₅	Ns
Agra	1604	736	600090062	585702444 (97.6%)	1925180	796802	79	199	0
Bpan	71	47	617043932	616226721 (99.8%)	111823772	110188975	3	4	4.0
Chom	522	82	534400685	521892830 (97.6%)	101522387	93884324	3	4	74.
Cruf	109329	75	288472566	4853570 (1.7%)	7933	4737	5734	13038	108
Cvic	197510	1	459231066	53150 (0.01%)	4480	3596	13675	24759	187
Cvom	129	41	706678423	704415709 (99.7%)	131076759	125523795	3	4	5.4
Lcup	8457	6	409159265	392766389 (95.9%)	71024062	67696999	3	4	49

Lser	4371	2222	565305478	513238679 (90.7%)	296078	137019	500	1197	0
Pazu	24	8	874252225	873869904 (99.9%)	164199054	157909393	3	4	8.8
Preg	192460	283	549932840	18344189 (3.3%)	13034	6995	8518	18763	104

Decisions

The Cvic, Preg, and Cruf genomes are too fragmented and this mitght not bee got for the WGA and the analysis of the non-coding regions. So I won't even mask them. Preg and Cruf might be ok for the gene family analysis (not Cvic), but this is something to be analyzed later.

3. Repeats and soft masking

I used RepeatModeler and RepeatMasker to get softmasked versions of all the genomes. Also, it provides tables summarizing the repeats found in each genome.

I did it all in Darwin, using Vanessa's user, because everything was already installed in there.

```
1 | # Genomes are here
2 | /home/cunha/Pedro/1-Genomes
```

3.1. RepeatModeler

Using Docker. One container for each genome, because we've had some porblems with containers desaparing before.


```

1  ##### General #####
2  # Get in
3  docker exec -it [name] /bin/bash
4  # Get out
5  Ctrl P+Q
6  # Check existing containers
7  docker ps
8
9  ##### Agra #####
10 # New container
11 docker run -it --rm dfam/tetools:latest # elastic_lehmann
12
13 # Copying genome to Docker
14 docker cp /home/cunha/Pedro/1-Genomes/Agra.fasta elastic_lehmann:/home
15
16 # Accessing the container
17 docker exec -it elastic_lehmann /bin/bash
18
19 # Database
20 BuildDatabase -name Agra_database Agra.fasta
21
22 # Running
23 RepeatModeler -database Agra_database -threads 30 -LTRStruct >Agra_log 2>Ag
24
25 # Copying the results
26 cd /home/cunha/Pedro/2-Masking/1-RepeatModeler
27
28 docker cp -a elastic_lehmann:/home/Agra_database-families.fa ./
29 docker cp -a elastic_lehmann:/home/Agra_database-families.stk ./
30 docker cp -a elastic_lehmann:/home/Agra_database-rmod.log ./
31 docker cp -a elastic_lehmann:/home/Agra_RepeatModeler.tar.gz ./ # this has a
32
33 ##### Bpan #####
34 # New container
35 docker run -it --rm dfam/tetools:latest # dreamy_hugle
36
37 # Copying genome to Docker
38 docker cp /home/cunha/Pedro/1-Genomes/Bpan.fasta dreamy_hugle:/home
39
40 # Accessing the container
41 docker exec -it dreamy_hugle /bin/bash
42
43 # Database
44 BuildDatabase -name Bpan_database Bpan.fasta
45
46 # Running
47 RepeatModeler -database Bpan_database -threads 30 -LTRStruct >Bpan_log 2>Bp
48
49 # Copying the results
50 cd /home/cunha/Pedro/2-Masking/1-RepeatModeler
51
52 docker cp -a dreamy_hugle:/home/Bpan_database-families.fa ./
53 docker cp -a dreamy_hugle:/home/Bpan_database-families.stk ./
54 docker cp -a dreamy_hugle:/home/Bpan_database-rmod.log ./
55 docker cp -a dreamy_hugle:/home/Bpan_RepeatModeler.tar.gz ./
56
57 ##### Chom #####
58 # Container - elastic_lehmann (same as Agra)
59
60 # Copying genome to Docker
61 docker cp /home/cunha/Pedro/1-Genomes/Chom.fasta elastic_lehmann:/home
62
63 # Accessing the container
64 docker exec -it elastic_lehmann /bin/bash
65
66 # Database
67 BuildDatabase -name Chom_database Chom.fasta
68
69 # Running
70 RepeatModeler -database Chom_database -threads 30 -LTRStruct >Chom_log 2>Ch
71
72 # Copying the results
73 cd /home/cunha/Pedro/2-Masking/1-RepeatModeler
74
75 docker cp -a elastic_lehmann:/home/Chom_database-families.fa ./
76 docker cp -a elastic_lehmann:/home/Chom_database-families.stk ./
77 docker cp -a elastic_lehmann:/home/Chom_database-rmod.log ./
78 docker cp -a elastic_lehmann:/home/Chom_RepeatModeler.tar.gz ./
79
80 ##### Cvom #####
81 # Container - dreamy_hugle (same as Bpan)
82
83 # Copying genome to Docker
84 docker cp /home/cunha/Pedro/1-Genomes/Cvom.fasta dreamy_hugle:/home
85
86 # Accessing the container
87 docker exec -it dreamy_hugle /bin/bash
88

```

```

89 # Database
90 BuildDatabase -name Cvom_database Cvom.fasta
91
92 # Running
93 RepeatModeler -database Cvom_database -threads 30 -LTRStruct >Cvom_log 2>Cv
94
95 # Copying the results
96 cd /home/cunha/Pedro/2-Masking/1-RepeatModeler
97
98 docker cp -a dreamy_hugle:/home/Cvom_database-families.fa ./
99 docker cp -a dreamy_hugle:/home/Cvom_database-families.stk ./
100 docker cp -a dreamy_hugle:/home/Cvom_database-rmod.log ./
101 docker cp -a dreamy_hugle:/home/Cvom_RepeatModeler.tar.gz ./
102
103 #### Lcup ####
104 # New container
105 docker run -it --rm dfam/tetools:latest # determined_ride
106
107 # Copying genome to Docker
108 docker cp /home/cunha/Pedro/1-Genomes/Lcup.fasta determined_ride:/home
109
110 # Accessing the container
111 docker exec -it determined_ride /bin/bash
112
113 # Database
114 BuildDatabase -name Lcup_database Lcup.fasta
115
116 # Running
117 RepeatModeler -database Lcup_database -threads 30 -LTRStruct >Lcup_log 2>Lc
118
119 # Copying the results
120 cd /home/cunha/Pedro/2-Masking/1-RepeatModeler
121
122 docker cp -a determined_ride:/home/Lcup_database-families.fa ./
123 docker cp -a determined_ride:/home/Lcup_database-families.stk ./
124 docker cp -a determined_ride:/home/Lcup_database-rmod.log ./
125 docker cp -a determined_ride:/home/Lcup_RepeatModeler.tar.gz ./
126
127 #### Lser ####
128 # New container
129 docker run -it --rm dfam/tetools:latest # happy_snyder
130
131 # Copying genome to Docker
132 docker cp /home/cunha/Pedro/1-Genomes/Lser.fasta happy_snyder:/home
133
134 # Accessing the container
135 docker exec -it happy_snyder /bin/bash
136
137 # Database
138 BuildDatabase -name Lser_database Lser.fasta
139
140 # Running
141 RepeatModeler -database Lser_database -threads 30 -LTRStruct >Lser_log 2>Ls
142
143 # Copying the results
144 cd /home/cunha/Pedro/2-Masking/1-RepeatModeler
145
146 docker cp -a happy_snyder:/home/Lser_database-families.fa ./
147 docker cp -a happy_snyder:/home/Lser_database-families.stk ./
148 docker cp -a happy_snyder:/home/Lser_database-rmod.log ./
149 docker cp -a happy_snyder:/home/Lser_RepeatModeler.tar.gz ./
150
151 # IMPORTANT NOT - THE LTR PIPELINE DID NOT WORK IN THIS AUTOMATED WAY
152
153 #### Pazu ####
154 # New container
155 docker run -it --rm dfam/tetools:latest # funny_payne
156
157 # Copying genome to Docker
158 docker cp /home/cunha/Pedro/1-Genomes/Pazu.fasta funny_payne:/home
159
160 # Accessing the container
161 docker exec -it funny_payne /bin/bash
162
163 # Database
164 BuildDatabase -name Pazu_database Pazu.fasta
165
166 # Running
167 RepeatModeler -database Pazu_database -threads 30 -LTRStruct >Pazu_log 2>Pa
168
169 # Copying the results
170 cd /home/cunha/Pedro/2-Masking/1-RepeatModeler
171
172 docker cp -a funny_payne:/home/Pazu_database-families.fa ./
173 docker cp -a funny_payne:/home/Pazu_database-families.stk ./
174 docker cp -a funny_payne:/home/Pazu_database-rmod.log ./
175 docker cp -a funny_payne:/home/Pazu_RepeatModeler.tar.gz ./
176

```



```

177 # Removing the containers (since I already got all the outputs)
178 docker stop elastic_lehmann # Agra & Chom
179 docker stop dreamy_hugle # Bpan & Cvom
180 docker stop determined_ride # Lcup
181 docker stop happy_snyder # Lser NOT YET
182 docker stop funny_payne # Pazu

```

I did all of this on Darwin, but I'll use RepeatMasker on Rosalind, so I'll move all the files there.

```

1 cd /home/cunha/Pedro/2-Masking/1-RepeatModeler # in Darwin
2
3 scp -P 2205 * pedro@143.107.244.181:/home/pedro/Non_Coding_Element_Evolution/

```

3.2. RepeatMasker

3.2.1. Installation

Pre requisites

```

1 # h5py python library
2 pip install h5py
3
4 # RMBlast v2.11.0
5 scp -P 2205 Mestrado_genomas/rmbblast-2.11.0+-x64-linux.tar.gz pedro@143.107.244.181:/home/pedro/Programs/
6
7 tar -xf rmbblast-2.11.0+-x64-linux.tar.gz
8
9 # Tandem Repeat Finder v4.09.1
10 scp -P 2205 Mestrado_genomas/TRF-4.09.1.tar.gz pedro
11 @143.107.244.181:/home/pedro/Programs
12
13 tar -xf TRF-4.09.1.tar.gz
14
15 cd TRF-4.09.1/
16 mkdir build
17 cd build
18 ../configure
19 make
20 sudo make install

```

RepeatMasker

```

1 # Distribution - RepeatMasker v4.1.2
2 scp -P 2205 Mestrado_genomas/RepeatMasker-4.1.5-pl.tar.gz pedro@143.107.244.181:/home/pedro/Programs/
3
4 tar -xf RepeatMasker-4.1.5-pl.tar.gz
5
6 cd home/pedro/Programs/RepeatMasker/
7
8 perl ./configure

```

About the additional libraries:

- The Dfam file (https://www.dfam.org/releases/Dfam_3.6/families/Dfam.h5.gz) is enormous and would take too long to download. Also, there already is a Dfam file in the RepeatMasker distribution, so it might work by itself (this is Dfam 3.3)
- The RepBase library requires a paid subscription to be downloaded, so we cannot use it.

3.2.2. Running

For this, I use the repeat libraries that I got using RepeatModeler, because they are species-specific.

```

1 | # I'm doing this on Rosalind
2 | cd /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker
3 |
4 | # repeatmasker.bash
5 | for GENOME in /home/pedro/Non_Coding_Element_Evolution/1-Genomes/*
6 | do
7 |     FASTA=${GENOME##*/home/pedro/Non_Coding_Element_Evolution/1-Genomes/}
8 |     DB=${FASTA%%.fasta}_database
9 |     SPP=${FASTA%%.fasta}
10 |
11 |     echo $SPP
12 |     echo "RepeatMasker -lib /home/pedro/Non_Coding_Element_Evolution/3-Masking/1-RepeatMaskerLib/$SPP
13 |
14 |     RepeatMasker -lib /home/pedro/Non_Coding_Element_Evolution/3-Masking/1-RepeatMaskerLib/$SPP
15 |
16 | done

```

With this command line, I had 6 outputs for each genome:

- Masked genome
- .align (alignments of all repeats to the genomes)
- .out (huge table with all information)
- .out.gff (outfile in gff format)
- .out.html (outfile in html format)
- .tbl (a smaller table that summarizes the .out table)

Organizing things:

```

1 | cd /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker
2 |
3 | # Moving some less important outputs to a subdirectory and compressing them
4 | mkdir Other_outputs
5 | mv *out* Other_outputs
6 | mv *align* Other_outputs
7 | mv *gz* Other_outputs
8 | mv *log* Other_outputs
9 | tar -czvf Other_outputs.tar.gz Other_outputs

```

LTR - Lser

When running RepeatModeler, the LTR pipeline failed for the Lser genome. So I did a separate analysis to solve this using LTR Harvest and LTR Retriever.

```

1 | docker run -it --rm dfam/tetools:latest # distracted_dewdney
2 |
3 | docker cp /home/pedro/Non_Coding_Element_Evolution/1-Genomes/Lser.fasta distracted_dewdney:/home/Lser.fasta
4 |
5 | docker exec -it distracted_dewdney /bin/bash
6 |
7 | # Got these command lines from the LTR retriever help guide
8 | # This is the LTR.sh script
9 | /opt/genometools/bin/gt suffixerator -db Lser.fasta -indexname Lser_db.fasta -l 100 -m 1000000
10 |
11 | /opt/genometools/bin/gt ltrharvest -index Lser_db.fasta -minlenltr 100 -maxlenltr 1000000
12 |
13 | /opt/LTR_retriever/LTR_retriever -genome Lser.fasta -inharvest Lser.harvest.fasta -out Lser.LTRlib.fasta

```

This worked! I have a fasta file with the identified LTR sequences from the Lser genome (and a bunch of other outputs). I copied everything back from the container.

```

1 | mkdir /home/pedro/Non_Coding_Element_Evolution/3-Masking/Lser_LTR
2 | cd /home/pedro/Non_Coding_Element_Evolution/3-Masking/Lser_LTR
3 |
4 | # Copy
5 | docker cp -a distracted_dewdney:/home/Lser.fasta.LTR.gff3 ./
6 | docker cp -a distracted_dewdney:/home/Lser.fasta.LTRlib.fasta ./
7 | docker cp -a distracted_dewdney:/home/Lser_LTR_outputs.tar.gz ./
8 | docker cp -a distracted_dewdney:/home/LTR.sh ./
9 |
10 | # Terminate the container
11 | docker stop distracted_dewdney

```

To check if it made any difference, I'll concatenate the Lser.fasta.LTRlib.fasta to the output I had from the first RepeatModeler run, then use this new, improved, library, for RepeatMasker.

```

1 | cd /home/pedro/Non_Coding_Element_Evolution/3-Masking/Lser_LTR
2 |
3 | cp /home/pedro/Non_Coding_Element_Evolution/3-Masking/1-RepeatModeler/Lser_d
4 |
5 | # New database
6 | cat Lser_database-families.fa Lser.fasta.LTRlib.fa >Lser_COMPLETE_database-fa

```

Finally, I do a new RepeatMasker run, to check if there was any difference from the previous trial:

```

1 | cd /home/pedro/Non_Coding_Element_Evolution/3-Masking/Lser_LTR
2 |
3 | RepeatMasker -lib Lser_COMPLETE_database-families.fa -dir ./ -pa 20 -a -xsmal

```

It was different (and better), so I'll use this new version of the masked genome.

```

1 | # Renaming and compressing the first run
2 | cd /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker/
3 |
4 | mkdir Lser_NO_LTR
5 |
6 | for i in Lser*; do mv $i Lser_NO_LTR${i##Lser}; done
7 | mv Lser_NO_LTR.* Lser_NO_LTR/
8 |
9 | cd Other_outputs/
10 | for i in Lser*; do mv $i Lser_NO_LTR${i##Lser}; done
11 | mv Lser_NO_LTR.* ../Lser_NO_LTR/
12 |
13 | cd ..
14 | tar -czvf Lser_NO_LTR.tar.gz Lser_NO_LTR
15 |
16 | # Moving things from the second run
17 | cd /home/pedro/Non_Coding_Element_Evolution/3-Masking/Lser_LTR
18 |
19 | mv Lser_COMPLETE_database-families.fa ../1-RepeatModeler/
20 |
21 | mv Lser.fasta.LTR.gff3 ../1-RepeatModeler/Lser_LTR/
22 | mv Lser.fasta.LTRlib.fa ../1-RepeatModeler/Lser_LTR/
23 | mv Lser_LTR_outputs.tar.gz ../1-RepeatModeler/Lser_LTR/
24 | mv LTR.sh ../1-RepeatModeler/Lser_LTR/
25 |
26 | mv Lser.fasta.masked Lser_masked.fasta
27 | mv Lser.fasta.tbl Lser_masked.tbl
28 | mv Lser_masked* ../2-RepeatMasker/
29 |
30 | mv * ../2-RepeatMasker/Other_outputs

```

3.3. Summarizing information

I organized the outputs using a python script I wrote and made some plots using an Rscript (I did this locally, but everything is on Rosalind now)

```

1 | # Everything is in here
2 | /home/pedro/Non_Coding_Element_Evolution/3-Masking/3-Organizing_information

```

3.4. Compressing unmasked genomes

Once it was finished. The unmasked genomes were compressed, and only the masked versions were kept for further analyses (this was made to take up less storage).

```

1 | /home/pedro/Non_Coding_Element_Evolution/1-Genomes.tar.gz

```

Note on the new genomes

We new genomes (Cmeg, Cput and Lexi) using the same pipeline used for the genomes I downloaded. From this moment on, I'm using all of them. I copied these masked genomes to my working directory and renamed them (to keep all files with the same patterns). I also renamed my own outputs to make everything simpler for me.

```

1 | cd /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker
2 |
3 | # Renaming my outputs
4 | for i in *masked; do mv $i ${i%.fasta.masked}_masked.fasta; done
5 |
6 | for i in *tbl; do mv $i ${i%.fasta.tbl}_masked.tbl; done
7 |
8 | # Renaming the outputs from the new genomes
9 | mv cput_N_genome_final.fa.masked Cput_masked.fasta
10 | mv cput_N_genome_final.fa.tbl Cput_masked.tbl
11 |
12 | mv cmeg_N_genome_final.fa.masked Cmeg_masked.fasta
13 | mv cmeg_N_genome_final.fa.tbl Cmeg_masked.tbl
14 |
15 | mv lexi_N_genome.fa.masked Lexi_masked.fasta
16 | mv lexi_N_genome.fa.tbl Lexi_masked.tbl

```

4. Whole genome alignment (WGA)

To find conserved non-coding elements, I had to create a WGA. To do so, I choose to use ProgressiveCactus. Briefly, it aligns the genomes based on a phylogeny. Pairwise alignments are made based on species relatedness and ancestral genomes are estimated on every node of the tree. Also, a final alignment of the genomes is provided.

General information on Cactus:

<https://github.com/ComparativeGenomicsToolkit/cactus>

4.1. Installation

Cactus

```

1 | # Precompiled binaries - Cactus v2.6.4
2 | wget https://github.com/ComparativeGenomicsToolkit/cactus/releases/download/v2.6.4/cactus-bin-v2.6.4.tar.gz
3 |
4 | tar -xf cactus-bin-v2.6.4.tar.gz
5 |
6 | cd /home/pedro/Programs/cactus-bin-v2.6.4/
7 |
8 | # Set up Python environment
9 | python3 -m pip install virtualenv
10 | virtualenv -p python3.8 cactus_env
11 |
12 | echo "export PATH=/home/pedro/Programs/cactus-bin-v2.6.4/bin:\$PATH" >> cactus_env/bin/activate
13 | echo "export PYTHONPATH=/home/pedro/Programs/cactus-bin-v2.6.4/lib:\$PYTHONPATH" >> cactus_env/bin/activate
14 |
15 | source cactus_env/bin/activate
16 |
17 | python3 -m pip install -U setuptools pip==23.2.1
18 | python3 -m pip install -U -r ./toil-requirement.txt
19 | python3 -m pip install -U .
20 |
21 | # Afterwards, that's how I can get in the environment:
22 | source /home/pedro/Programs/cactus-bin-v2.6.4/cactus_env/bin/activate
23 | # This is how to get out of it:
24 | deactivate

```

HAL

I installed HAL tools so I could turn the Cactus HAL output (a reference free alignment) into other formats that can be used in other analyses.

```

1  # HDF5 1.10.1 with C++ API enabled
2  mkdir /home/pedro/Programs/hdf5
3  wget http://www.hdfgroup.org/ftp/HDF5/releases/hdf5-1.10/hdf5-1.10.1/src/hdf5-1.10.1.tar.gz
4  tar xzf hdf5-1.10.1.tar.gz
5  cd hdf5-1.10.1
6  ./configure --enable-cxx --prefix /home/pedro/Programs/hdf5
7  make && make install
8
9  export PATH=/home/pedro/Programs/hdf5/bin:${PATH}
10 export h5prefix=prefix=/home/pedro/Programs/hdf5
11
12 #sonLib
13 git clone https://github.com/ComparativeGenomicsToolkit/sonLib.git
14 pushd sonLib && make sonLibRootDir=/home/pedro/Programs/sonLib && popd
15
16 pip install sonLib
17
18 # Gitclone - HAL v2.2
19 cd /home/pedro/Programs
20 git clone https://github.com/ComparativeGenomicsToolkit/hal.git
21 cd hal
22 make
23 export PATH=/home/pedro/Programs/hal/bin:${PATH}
24 export PYTHONPATH=/home/pedro/Programs:${PYTHONPATH}
25 export PYTHONPATH=/home/pedro/Programs/hal:${PYTHONPATH}

```

cwltool

I had to install this one, because a message came up when I was performing the alignments (not a fatal error, but it's better to get things right)

```

1  # I need to this inside the environment
2  source /home/pedro/Programs/cactus-bin-v2.6.4/cactus_env/bin/activate
3
4  pip install cwltool

```

Normally, this would suffice. But I got a message suggesting that I used a different version to match other dependencies. I'm assuming this varies depending on the situation, but here's the information anyway:

```

1  pip install cwltool==3.1.20211107152837

```

4.2. Tree

Cactus requires a phylogenetic tree to perform the WGA. The tree Gisele and I used doesn't include all the species I'm working with for this analysis, nor does the Dimensions tree (which also has the problem of having more than 1 branch for some species). So I'll use a new tree using the BUSCO outputs. I'm doing this with Vanessa and Letícia. All the information is [here](#).

After having the tree, I slightly edited it. I removed unwanted tips (the species that have no genomic information), and rewrote some names (some had lowercase letters)

```

1  # cactus_tree.R
2  library(ape)
3  library(phytools)
4
5  setwd("C:/Users/Pedro/OneDrive/Área de Trabalho/Cactus")
6
7  # Read tree
8  tree <- read.tree("Calliphoridae_tree_root.tre")
9  plotTree(tree) # check
10
11 # Remove branches that won't be used
12 cactus <- drop.tip(tree, c("clop", "loch", "cbez", "calb"))
13 plotTree(cactus) # check
14
15 # Rename branches with lowercase letters
16 cactus$tip.label[1] <- "Cmeg"
17 cactus$tip.label[2] <- "Cput"
18 cactus$tip.label[7] <- "Lexi"
19 plotTree(cactus) # check
20
21 # Newick
22 write.tree(cactus, "Cactus_tree.tre")

```

4.3. Running Cactus

Cactus file

Cactus relies on a specific kind of file that simultaneously contains the tree and the pathway to each genome. Also, I can choose some genomes to be marked as having reference quality: I chose all the ones that had an $L_{50}=3$ (see QUAST results)

The file is here:

```
1 cd /home/pedro/Non_Coding_Element_Evolution/4-ProgressiveCactus_WGA/
2
3 # Cactus_file.txt
4 (((Cmeg:0.0184171866,Cput:0.0451563064):0.0294538378,Chom:0.059936427):0.0000000000,
5
6 Agra /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker/Agra_r
7 *Bpan /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker/Bpan_r
8 *Cvom /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker/Cvom_r
9 Cmeg /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker/Cmeg_r
10 Cput /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker/Cput_r
11 *Chom /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker/Chom_r
12 Lcup /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker/Lcup_r
13 Lser /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker/Lser_r
14 Lexi /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker/Lexi_r
15 *Pazu /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker/Pazu_r
```

Alignment

```
1 cd /home/pedro/Non_Coding_Element_Evolution/4-ProgressiveCactus_WGA
2
3 source /home/pedro/Programs/cactus-bin-v2.6.4/cactus_env/bin/activate
4
5 # Running
6 cactus NEW_WGA Cactus_file.txt Calliphoridae_WGA.hal --maxCores 20 --logCrit:
7
8 # After completion
9 deactivate
10
11 # Turn it to MAF format (one MAF file for each scaffold in the genome)
12 export PATH=/home/pedro/Programs/cactus-bin-v2.6.4/bin:${PATH}
13
14 hal2mafMP.py --numProc 5 --splitBySequence --refGenome Chom --noAncestors --
15
16 # Better names
17 for i in *maf; do mv $i ${i#Calliphoridae_WGA_}; done
18
19 # Remove second line of the files, because it has a tree with the ancestors
20 for i in *maf; do sed -i -e 2d $i; done
21
22 # Keep the files elsewhere
23 mkdir MAF_files_Chom_ref
24 mv *.maf MAF_files_Chom_ref/
```

5. PhastCons

I'll use the PHAST software (more specifically, PhastCons) to find conserved elements in the WGA of Calliphoridae species.

IMPORTANT REFERENCES:

<http://compugen.cshl.edu/phast/resources.php>

<http://compugen.cshl.edu/phast/phastCons-HOWTO.html#paper>

https://github.com/tsackton/ratite-genomics/tree/master/04_wga/02_ce_id

5.1. Installation

```
1 # Clapack v3.2.1 (required for PHAST)
2 cd /home/pedro/Programs
3
4 wget http://www.netlib.org/clapack/clapack.tgz
5 tar -xf clapack.tgz
6
7 cd CLAPACK-3.2.1/
8 cp make.inc.example make.inc && make f2clib && make blaslib && make lib
9
10 # PHAST v1.6
11 scp -P 2205 phast-1.6.tar.gz pedro@143.107.244.181:/home/pedro/Programs
12 tar -xf phast-1.6.tar.gz
13
14 cd phast-1.6/src
15 make CLAPACKPATH=/home/pedro/Programs/CLAPACK-3.2.1
```

5.2. Running

In the PHAST folder I have the WGA in MAF format (with a *Co. hominivorax* reference and without the ancestral genome estimates - see Cactus section), the *Co. hominivorax* genome gff file, and the tree I used for the WGA.


```

1 cd /home/pedro/Non_Coding_Element_Evolution/5-PHAST
2
3 # Now, to make the conserved and non-conserved models and predict the element
4
5 cd /home/pedro/Non_Coding_Element_Evolution/5-Phast/1-Chom_per_Scaffold
6
7 cp /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-RepeatMasker/Chom_masked.fasta
8
9 # Making the sequence names equal to the MAF files
10 sed -i 's/Chom_/Chom.Chom_/g' Chom_masked.fasta
11
12 # Split fasta (found it here: https://gist.github.com/astatham/621901)
13 cat Chom_masked.fasta | awk '{if (substr($0, 1, 1)==">") {filename=(substr($0, 2, 100));}}' > Chom_masked_split.fasta
14
15 # I did this, because I changed the sequence headers earlier, but it's best to keep them as is
16 for i in *.fa; do mv $i ${i#Chom.};done
17
18 # Now, do the splitting - chunks.sh
19 cd /home/pedro/Non_Coding_Element_Evolution/5-Phast/2-MAF_chunks
20
21 for MAF in /home/pedro/Non_Coding_Element_Evolution/4-ProgressiveCactus_WGA/MAF/*; do
22
23 do
24     SCAFF=${MAF##*/home/pedro/Non_Coding_Element_Evolution/4-ProgressiveCactus_WGA/MAF/}
25     SCAFF=${SCAFF%%.maf}
26
27     msa_split $MAF --in-format MAF --refseq /home/pedro/Non_Coding_Element_Evolution/4-ProgressiveCactus_WGA/MAF/$SCAFF --out-dir /home/pedro/Non_Coding_Element_Evolution/5-Phast/2-MAF_chunks/$SCAFF
28 done
29
30 # GC content
31 # Get the average GC content of the WGA to use it in the model prediction
32 cd /home/pedro/Non_Coding_Element_Evolution/5-Phast/3-Models
33
34 mkdir basecomp
35 cd basecomp
36
37 export PATH=/home/pedro/Programs/hal/bin:${PATH}
38
39 for spp in $(halStats --genomes /home/pedro/Non_Coding_Element_Evolution/4-ProgressiveCactus_WGA/MAF/*); do
40
41 #halStats --basecomp -> the output is fraction_of_As fraction_of_Gs fraction_of_Cs fraction_of_Ts
42
43 # I don't need the information on the ancestral genomes, I'll put it away for later
44 mkdir Ancestral_genomes
45 mv Anc* Ancestral_genomes/
46
47 # Average CG content in my genomes
48 # I need columns 2 and 3, because they are G and C
49 cat *basecomp | awk '{SUM+= $2; SUM+= $3; print SUM/11}' | tail -n 1 > GC # Average GC content
50
51 # Model prediction
52 cd /home/pedro/Non_Coding_Element_Evolution/5-Phast/3-Models
53
54 mkdir Trees
55
56 # model_predict.sh
57 for MAF in /home/pedro/Non_Coding_Element_Evolution/5-Phast/2-MAF_chunks/*; do
58
59 do
60     NAME=${MAF##*/home/pedro/Non_Coding_Element_Evolution/5-Phast/2-MAF_chunks/}
61     NAME=${NAME%%.maf}
62
63     phyloFit --tree /home/pedro/Non_Coding_Element_Evolution/4-ProgressiveCactus_WGA/MAF/$NAME --msa-format SS $MAF --out-dir /home/pedro/Non_Coding_Element_Evolution/5-Phast/3-Models/Trees/$NAME
64
65     phastCons --gc 0.26 --estimate-trees Trees/$NAME --msa-format SS $MAF --out-dir /home/pedro/Non_Coding_Element_Evolution/5-Phast/3-Models/Trees/$NAME
66 done
67
68 # Organizing the initial models
69 cd /home/pedro/Non_Coding_Element_Evolution/5-Phast/3-Models
70 mkdir Initial_models
71 mv init.* Initial_models
72
73 # General organization
74 mv basecomp 1-Basecomp
75 mv Initial_models 2-Initial_models
76 mv Trees 3-Trees
77
78 # I need to filter the models, leaving only the ones with all species (because of the way the models are named)
79 cd 3-Trees
80
81 mkdir not_all_spp
82
83 for MOD in *mod
84 do
85     TREE=$(tail -n 1 $MOD)
86     SPP='Cmeg'*'Cput'*'Chom'*'Pazu'*'Lser'*'Lcup'*'Lexi'*'Cvom'*'Agra'*'Bpar'
87
88     if [[ $TREE != *$SPP* ]]

```

```

89     then
90         mv $MOD not_all_spp/
91     fi
92 donefor MOD in *mod
93 do
94     TREE=$(tail -n 1 $MOD)
95     SPP='Cmeg'*'Cput'*'Chom'*'Pazu'*'Lser'*'Lcup'*'Lexi'*'Cvom'*'Agra'*'Bpar
96
97     if [[ $TREE != *$SPP* ]]
98     then
99         mv $MOD not_all_spp/
100     fi
101 done
102
103 # Now I average the models to do a final run
104 cd /home/pedro/Non_Coding_Element_Evolution/5-Phast/3-Models
105 mkdir 4-Average_models
106 cd 4-Average_models
107
108 ls ../3-Trees/*.cons.mod > cons.txt
109 phyloBoot --read-mods '*cons.txt' --output-average ave_cons.mod
110
111 ls ../3-Trees/*.noncons.mod > noncons.txt
112 phyloBoot --read-mods '*noncons.txt' --output-average ave_noncons.mod
113
114 # Naming the ancestral nodes in the model (will need it afterwards)
115 tree_doctor --name-ancestors ave_cons.mod > ave_cons_named.mod
116 tree_doctor --name-ancestors ave_noncons.mod > ave_noncons_named.mod
117
118 # Final run with the averaged models
119 cd /home/pedro/Non_Coding_Element_Evolution/5-Phast/3-Models
120 mkdir 5-Final_predictions
121 cd 5-Final_predictions
122
123 mkdir Elements Scores
124
125 for MAF in /home/pedro/Non_Coding_Element_Evolution/5-Phast/2-MAF_chunks/*ss
126 do
127     NAME=${MAF##/home/pedro/Non_Coding_Element_Evolution/5-Phast/2-MAF_chunk
128     NAME=${NAME%.ss}
129
130     phastCons --gc 0.26 --most-conserved Elements/$NAME.bed --score --msa-fc
131 done
132
133 cd Elements
134
135 cat Chom*.bed | sort -k1,1 -k2,2n > most-conserved.bed
136
137 # 818,864 elements found

```

5.3. Element filtering

I used bedtools to filter my conserved elements in respect to the whole genome

5.3.1. Installing bedtools and bedops

```

1 # Now i need bedtools to get the elements that are not within genes
2 # bedtools v2.30.0
3 scp -P 2205 bedtools-2.30.0.tar.gz pedro@143.107.244.181:/home/pedro/Programs
4 tar -xvf bedtools-2.30.0.tar.gz
5
6 cd bedtools2
7 make
8
9 # BEDOPS v2.4.41
10 wget https://github.com/bedops/bedops/releases/download/v2.4.41/bedops_linux_
11
12 cd bedops_v2.4.41

```

5.3.2. Filtering

```

1 # Moving on
2 mkdir /home/pedro/Non_Coding_Element_Evolution/5-Phast/4-Conserved_elements
3 cd /home/pedro/Non_Coding_Element_Evolution/5-Phast/4-Conserved_elements
4
5 mv ../3-Models/5-Final_predictions/Elements/most-conserved.bed ./all_elements
6
7 # I also merged the element predictions that were 5bp or less away from each
8 bedtools merge -i all_elements.bed -d 5 -c 4,5,6 -o distinct > merged_elements
9
10 # Filtering elements with at least 250bp
11 awk '($3 - $2) >= 250' merged_elements.bed > 250bp_elements.bed # 44,253 (250
12
13 # Now, to keep only the elements that are non-exonic, I need to use the Co. l
14 mkdir Chom_ref
15 cd Chom_ref
16
17 cp /home/blowflies/Gene_families/00-Genomic_data/05-GFF/Chom_GFF.gtf .
18
19 sed -i 's/Chr/Chom_Chrg/' Chom_GFF.gtf # just to make the gff equal to the p
20
21 # Getting just the exon entries
22 grep "exon" Chom_GFF.gtf > Chom_exons.gtf
23
24 # Convert to bed
25 gtf2bed < Chom_exons.gtf > Chom_exons.bed
26
27 # When I tried to use bedops on the entire genome gff, I was warned that the
28 sed -r 's/\\;\\.*/\\/' Chom_exons.bed > Chom_exons_corrected.bed
29
30 # Get non-exonic elements
31 cd ..
32 bedtools intersect -a 250bp_elements.bed -b Chom_ref/Chom_exons_corrected.bed

```

All	Merged	250bp	Non-exonic
818,864	797,186	44,253	20,785

6. Element processing

6.1. Sequences in the other species

Now that I identified the CNEEs with the Co. hominivorax reference, I need to find them in the WGA for the other species. After that I can align them and make the following analyses.

```

1 # Finding the CNEEs in the other species
2 mkdir /home/pedro/Non_Coding_Element_Evolution/6-Element_processing
3 cd /home/pedro/Non_Coding_Element_Evolution/6-Element_processing
4
5 # hal2bed.sh
6 for SP in Cmeg Cput Pazu Lser Lcup Lexi Cvom Agra Bpan
7 do
8     halLiftover /home/pedro/Non_Coding_Element_Evolution/4-ProgressiveCactus
9 done
10
11 for i in *bed; do cat $i | sort -k1,1 -k2,2n > ${i%.bed}_sorted.bed; done

```

Some Co. hominivorax CNEEs had more than one correlate in other species (a few bases apart), so merging using bedtools didn't work well. So I used a costum python script to do the merging (bed_file_organizer.py).

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3
4  '''
5  PROGRAM: bed_file_organizer.py
6  '''
7  import glob
8  import sys
9  import getopt
10 #-----
11 def main():
12     pwd = arg()
13
14     if pwd == None:
15         usage()
16
17
18     bed_files = glob.glob(f'{pwd}/*bed') # getting all files
19
20     for bed in bed_files:
21
22         start = None
23         end = None
24         ref = None
25         merge={}
26
27         file = open(bed, 'r')
28         out = open(f'{bed}_merged.bed', 'w')
29
30         for l in file:
31             info = l.split()
32
33             chrom = info[0]
34             start = int(info[1])
35             end = int(info[2])
36             ref = info[3]
37             n = info[4]
38             strand = info[5]
39
40             name=chrom+ref
41
42             if name not in merge:
43                 merge[name] = [chrom, [start], [end], ref, n, strand]
44             else:
45                 merge[name][1]+=[start]
46                 merge[name][2]+=[end]
47
48             for i in merge:
49                 s = min(merge[i][1])
50                 e = max(merge[i][2])
51
52                 out.write(f'{merge[i][0]}\t{s}\t{e}\t{merge[i][3]}\t{merge[i][4]}
53
54             file.close()
55             out.close()
56
57 #-----
58 def arg():
59     argv = sys.argv[1:]
60     try:
61         opts,args = getopt.getopt(argv, 'p:')
62     except:
63         usage()
64
65     arg_1 = None
66
67     for opt, arg in opts:
68         if opt in ['-p']:
69             arg_1 = arg
70
71     return arg_1
72 #-----
73 def usage():
74     sys.exit('Usage:\n'
75             '\t-p: files directory pathway\n')
76 #-----
77 if __name__ == '__main__':
78     main()

```

```

1 cd /home/pedro/Non_Coding_Element_Evolution/6-Element_processing
2
3 python3 bed_file_organizer.py -p .
4
5 # Fixing the output names
6 for i in *merged*; do mv $i ${i%.bed_merged.bed}_merged.bed; done
7
8 # Organizing
9 mkdir Unsorted Sorted Merged
10 mv *merged.bed Merged/
11 mv *sorted.bed Sorted/
12 mv *.bed Unsorted/
13
14 mkdir 1-Bed_other_spp
15 mv * 1-Bed_other_spp

```

Species	# CNEEs
Agra	16,562
Bpan	16,628
Cmeg	21,285
Cput	19,918
Cvom	16,599
Lcup	14,913
Lexi	21,710
Lser	15,748
Pazu	18,726

Now that I have the bed files for each species, I can generate fasta files:

```

1 mkdir /home/pedro/Non_Coding_Element_Evolution/6-Element_processing/2-Fasta_files
2 cd /home/pedro/Non_Coding_Element_Evolution/6-Element_processing/2-Fasta_files
3
4 # Chom (because this is the reference, I'm using the final CNEEs I identified)
5 bedtools getfasta -fi /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-Reference.bed -bed /home/pedro/Non_Coding_Element_Evolution/6-Element_processing/3-Bed_files/Chom.bed -fo Chom.fasta
6
7 # Other species
8 for SP in Cmeg Cput Pazu Lser Lcup Lexi Cvom Agra Bpan; do bedtools getfasta -fi /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-Reference.bed -bed /home/pedro/Non_Coding_Element_Evolution/6-Element_processing/3-Bed_files/$SP.bed -fo $SP.fasta; done
9
10 # I'll put each species name in the fasta headers to make it easier for when I need to filter
11 for SP in Agra Bpan Chom Cmeg Cput Cvom Lcup Lexi Lser Pazu; do sed -i 's/:/:$SP/' $SP.fasta; done
12 sed -i 's/,/_/g' *fasta # the comma would be a problem afterwards
13
14 # Now, I need one file for each CNEE
15 mkdir 1-Species_files 2-CNEE_files
16 mv *fasta 1-Species_files
17
18 cd 2-CNEE_files
19 cut -f4 /home/pedro/Non_Coding_Element_Evolution/5-Phast/4-Conserved_elements.bed > CNEE_list.txt
20
21 # Fixing things that might cause trouble with grep
22 sed -i 's/$/:/' CNEE_list.txt
23 sed -i 's/,/_/g' CNEE_list.txt
24 sed -i 's/^Chom/INICHom/' CNEE_list.txt
25
26 # Also fixing the species files
27 sed -i 's/>Chom/>INICHom/' ../1-Species_files/*fasta
28
29 # Now, the CNEE files
30 while read line; do seqkit grep -r -p "$line" ../1-Species_files/*fasta >> $line.fasta; done < CNEE_list.txt
31
32 # Making the fasta files a bit better
33 sed -i 's/:/:@/2/' *fasta
34 sed -i 's/>.*:./>/' *fasta
35
36 # Also changing their name
37 for i in *fasta; do mv $i ${i#INI}; done

```

6.2. Filtering fasta files

Now everything worked! I have to filter these files so I can align them, and these are my criteria:

- No duplicates - just 1 sequence per species
- Sequences have to be $<2x$ and $>0.5x$ the *Co. hominivorax* reference
- All species have to be present

```

1  cd /home/pedro/Non_Coding_Element_Evolution/6-Element_processing/2-Fasta_files
2  mkdir Dupes No_dupes
3
4  # no_dupes.sh - selecting the ones with no duplicate species
5  for i in *.fasta
6  do
7      a=$(grep -c ">" $i) # number of sequences
8      b=$(grep ">" $i | sed 's/> //' | sed -r 's/@.*//' | uniq | wc -l) # number of species
9
10     if [ $a == $b ] # avoid duplicates (seq = spp)
11     then
12         mv $i No_dupes
13     fi
14 done
15
16 mv *fasta Dupes # just what is left
17
18 # size_filter.py - filtering for size (remove sequences that are more than 2x
19 cd No_dupes
20 for i in *fasta; do python3 ../size_filter.py $i; done
21
22 cd ..
23 mkdir Size_filtered
24 mv No_dupes/*filt*fasta Size_filtered

```

Just for the record:

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3
4  '''
5  PROGRAM: size_filter.py
6  '''
7
8  from Bio import SeqIO
9  import sys
10
11 def main(input_file):
12     reference_length = None
13
14     # Primeira iteração para encontrar o tamanho da sequência de referência
15     for record in SeqIO.parse(input_file, "fasta"):
16         if "Chom" in record.id:
17             reference_length = len(record.seq)
18             break
19
20     if reference_length is None:
21         print("Sequência de referência não encontrada em", input_file)
22         return
23
24     # Segunda iteração para escrever as sequências que atendem ao critério
25     output_file = input_file.replace(".fasta", "_filtered.fasta")
26     with open(output_file, "w") as output:
27         for record in SeqIO.parse(input_file, "fasta"):
28             if reference_length / 2 <= len(record.seq) <= reference_length * 2:
29                 SeqIO.write(record, output, "fasta")
30
31 if __name__ == "__main__":
32     if len(sys.argv) != 2:
33         print("Uso: python3 size_filter.py arquivo.fasta")
34         sys.exit(1)
35
36     input_file = sys.argv[1]
37     main(input_file)

```

Back to the filtering process:

```

1  cd /home/pedro/Non_Coding_Element_Evolution/6-Element_processing/2-Fasta_files
2  mkdir /home/pedro/Non_Coding_Element_Evolution/6-Element_processing/2-Fasta_files
3
4  # n_seq_filter.sh - filtering for number of sequences
5  for i in *filt*.fasta
6  do
7      a=$(grep -c ">" $i) # number of sequences
8
9      if [ $a -eq 10 ] # all species
10     then
11         cp $i /home/pedro/Non_Coding_Element_Evolution/6-Element_processing/2-Fasta_files
12     fi
13 done
14 # 3851

```

6.3. Alignments

I need to align each CNEE before I test for the acceleration.

```

1  mkdir /home/pedro/Non_Coding_Element_Evolution/6-Element_processing/3-Alignments
2  cd /home/pedro/Non_Coding_Element_Evolution/6-Element_processing/3-Alignments
3
4  # mafft_aln.sh
5  for FASTA in /home/pedro/Non_Coding_Element_Evolution/6-Element_processing/2-Fasta_files
6  do
7      FILE=${FASTA##*/} /home/pedro/Non_Coding_Element_Evolution/6-Element_processing/3-Alignments
8
9      mafft-ginsi --thread 10 --maxiterate 1000 --adjustdirectionaccurately --
10 done
11
12 sed -i 's/_R_/' * # because of the --adjustdirectionaccurately option

```

7. PhyloACC

Now I'll use PhyloACC to find elements that are accelerated in *Co. hominivorax* and/or *Pr. azurea*. For this I need the aligned CNEEs, a neutral model (I'll use the averaged non-conserved model for the WGA).

7.1. Installing

```

1  # PhyloACC v2.2.0
2  conda create -n phylo -c bioconda
3  conda activate phylo
4  conda install phyloacc
5
6  # astral v5.7.1
7  cd /home/pedro/Programs
8  wget https://github.com/smirarab/ASTRAL/archive/refs/tags/v5.7.1.tar.gz
9  tar -xvzf v5.7.1.tar.gz
10 cd ASTRAL-5.7.1
11 bash make.sh

```

7.2. Running

7.2.1. Coalescence tree

PhyloACC requires a tree in coalescent units, and I'll have to estimate that (and it has to have the same topology as the tree in the Phast model). I estimated it using the entire phylogeny (the I used for the WGA)

```

1  # Doing it on Darwin
2  cd /home/martins/Programs/ASTRAL-5.7.1
3
4  scp -P 2205 pedro@143.107.244.181:/home/pedro/Non_Coding_Element_Evolution/4-
5
6  java -jar astral.5.7.1.jar -i Calliphoridae_tree.tre -o Calliphoridae_astral
7
8  # The root in the Astral tree is non-significative, so I rerooted in in R
9  astral <- read.tree(text="(cbez,(cmeg,((calb,cput)0.67:0.2876820724517809,(l
10
11  a <- reroot(astral,21)
12
13  astral_new <- drop.tip(a, c("clop","loch","cbez","calb"))
14
15  write.tree(astral_new)
16
17  "(((Cmeg:NaN,Cput:NaN)0.67:0.2876820725,Chom:NaN)0.67:0.2876820725,Pazu:NaN)
18
19  # I had to move the branches (respecting the topology), so the species were :
20
21  "(((Cmeg:NaN,Cput:NaN)0.67:0.2876820725,Chom:NaN)0.67:0.2876820725,Pazu:NaN)
22
23  # Back to Darwin
24  cd /home/martins/PhyloACC
25  nano Astral_tree.tre
26  (((Cmeg:NaN,Cput:NaN)0.67:0.2876820725,Chom:NaN)0.67:0.2876820725,Pazu:NaN)
27
28  sed -i 's/:NaN//g' Astral_tree.tre
29
30  # Also, I removed the ancestor "names" (0.67), because this was also an issue
31  sed -i 's/0\.67:/:/g' Astral_tree.tre
32
33  (((Cmeg,(Calb,Cput)Calb-Cput:0.2876820725)Cmeg-Calb:0.2876820725,Chom)Cmeg-C

```

7.2.2. Other important files

```

1  cd /home/martins/PhyloACC
2
3  mkdir Aln
4
5  # Bringing alignments
6  scp -P 2205 pedro@143.107.244.181:/home/pedro/Non_Coding_Element_Evolution/6-
7
8  sed -i 's/@.*//' *
9
10 # Getting the model file
11 scp -P 2205 pedro@143.107.244.181:/home/pedro/Non_Coding_Element_Evolution/5-

```

7.2.3. Running


```

1  # PhyloACC
2  cd /home/martins/PhyloACC
3  conda activate phylo
4
5  # Preparing the data
6  phyloacc.py -d /home/martins/PhyloACC/Aln -m ave_noncons_named.mod -l Astral
7
8  # Running - Option 1
9  # snakemake -p -s /home/martins/PhyloACC/PhyloBoth/phyloacc-job-files/snakem
10
11 # Running - Option 2
12 max_attempts=100 # Maximum number of attempts
13
14 for numero in {1..78}; do # Loop for batches 1-78
15     echo "Tentando análise para número $numero:"
16
17     # Loop to try many times (if needed)
18     for ((i=1; i<=max_attempts; i++)); do
19         # Actual command line
20         PhyloAcc-GT PhyloBoth/phyloacc-job-files/cfgs/${numero}-gt.cfg > Phy
21
22         # Did the command work? (0 = success)
23         if [ $? -eq 0 ]; then
24             echo "Análise bem-sucedida para número $numero!"
25             break # Break loop if it works
26         else
27             echo "Erro de segmentação. Tentando novamente para número $numero"
28             sleep 1 # Wait a second before trying again (optional)
29         fi
30     done
31 done
32
33 # Post-processing the results
34 phyloacc_post.py -i PhyloBoth

```

Note:

The *phyloacc.py* command creates batches of alignments to allow the next steps and the snakemake files to run each batch. But it didn't work as expected, as some batches resulted in segmentation faults, which not only halted their execution, but the execution of all batches.

The authors suggested that I run each batch individually and then do the final step afterwards. Just as a reference, I left the original snakemake command line as a comment (Running - Option 1), but what was actually done is the the second one (Running - Option 2). Renan helped in the loop!

7.2.4. Back to Rosalind

```

1  scp -r -P 2205 /home/martins/PhyloACC pedro@143.107.244.181:/home/pedro/Non_C
2
3  mv PhyloACC/ 7-PhyloAcc/
4
5  # Some organization
6  cd 7-PhyloAcc
7  mkdir Slurm_files
8  mv phylo*.* Slurm_files

```

I simplified the *elem_lik.txt* result file (*/home/pedro/Non_Coding_Element_Evolution/7-PhyloAcc/PhyloBoth/results/elem_lik.txt*). But I did this locally, using Excel. I uploaded the files this directory: *home/pedro/Non_Coding_Element_Evolution/7-PhyloAcc/Result_files* . Briefly, I filtered all elements for which M1 was the best fit. Then, I split this filtered file into smaller ones: accelerated in *Co. hominivorax*; accelerated in *Pr. azurea*; accelerated in both; accelerated in none.

8. Closest genes

Now that I identified the elements that have accelerated evolution, I wanted to find the genes whose start is closest to these elements.

8.1. Bed files

I filtered the original bed file with all the elements to leav only the interesting ones.

```

1 cd /home/pedro/Non_Coding_Element_Evolution
2 mkdir 8-Closest_genes
3 cd 8-Closest_genes
4 mkdir 1-Filtered_beds
5
6 # Bed file with all elements found
7 cp ~/Non_Coding_Element_Evolution/5-Phast/4-Conserved_elements/merged_elements.bed
8
9 # Lists of accelerated elements
10 cut -f2 ~/Non_Coding_Element_Evolution/7-PhyloAcc/Result_files/M1_both.txt > M1_both.txt
11 cut -f2 ~/Non_Coding_Element_Evolution/7-PhyloAcc/Result_files/M1_all.txt > M1_all.txt
12 cut -f2 ~/Non_Coding_Element_Evolution/7-PhyloAcc/Result_files/M1_Chom.txt > M1_Chom.txt
13 cut -f2 ~/Non_Coding_Element_Evolution/7-PhyloAcc/Result_files/M1_Pazu.txt > M1_Pazu.txt
14 cut -f2 ~/Non_Coding_Element_Evolution/7-PhyloAcc/Result_files/M1_none.txt > M1_none.txt
15
16 # Some editions
17 sed -i -e ld *txt
18 sed -i 's/_filtered_aln/' *txt
19
20 sed -i 's/,/_/g' merged_elements.bed
21
22 # New bed files
23 for i in *txt
24 do
25     while read line
26     do
27         grep -w $line merged_elements.bed >> ${i%.txt}.bed
28         done < $i
29     done
30
31 for i in Elem*bed; do bedtools sort -i $i > ${i%.bed}_sorted.bed; done

```

8.2. Transcription start sites

First, I need the entire *Co. hominivorax* gff, so I can extract the transcription start sites.

```

1 mkdir /home/pedro/Non_Coding_Element_Evolution/8-Closest_genes/2-Transcription_start_sites
2 cd /home/pedro/Non_Coding_Element_Evolution/8-Closest_genes/2-Transcription_start_sites
3
4 # I extracted the first CDS from each gene, so I knew what is the transcript_id
5
6 cp /home/pedro/Non_Coding_Element_Evolution/5-Phast/4-Conserved_elements/Chom_exons_corrected.bed
7
8 bedtools sort -i Chom_exons_corrected.bed > Chom_exons_sorted.bed
9
10 sed -i 's/"/"/g' Chom_exons_sorted.bed
11 sed -i 's/transcript_id/' Chom_exons_sorted.bed
12
13 bedtools groupby -i Chom_exons_sorted.bed -g 4 -c 1,2,3,4,5,6 -o first > Chom_exons_grouped.bed
14
15 # Now I want just the first site
16 awk '{print $2, $3, $3 + 1, $5, $6, $7}' Chom_exons_grouped.bed > Chom_transcription_start.bed
17
18 # For some reason, the tabs were lost
19 sed -i "s/ /\t/g" Chom_transcription_start.bed

```

8.3. Closest genes

```

1 mkdir /home/pedro/Non_Coding_Element_Evolution/8-Closest_genes/3-Closest_genes
2 cd /home/pedro/Non_Coding_Element_Evolution/8-Closest_genes/3-Closest_genes
3
4 for i in ../1-Filtered_beds/Elem*sorted.bed; do bedtools closest -d -a $i -b
5
6 for i in *; do mv $i ${i%.sorted.bed}_closest.bed; done
7
8 # Filtering the closest gene and the distance to the closest CNEE
9 for i in *bed
10 do
11     echo -e "CNEE\tGene\tDistance" > ${i%.bed}_list.txt
12     cut -f4,10,13 $i >> ${i%.bed}_list.txt
13 done

```

9. Transcription factors

I checked whether the accelerated CNEEs were enriched for transcription factor binding motifs. And then, I checked for GO terms associated to genes that are controlled by TFs binding to these motifs.

9.1. Getting the CNEEs

```
1 | cd /home/pedro/Non_Coding_Element_Evolution/8-Closest_genes/1-Filtered_beds
2 |
3 | bedtools getfasta -fi /home/pedro/Non_Coding_Element_Evolution/3-Masking/2-R
```

9.2. TF binding motifs

MEME Suite - AME: <https://meme-suite.org/meme/tools/ame>

❗ Image Not Showing

Possible Reasons

- The image was uploaded to a note which you don't have access to
- The note which the image was originally uploaded to has been deleted

[Learn More →](#)

9.3. GO enrichment for genes with the TF binding motifs

MEME Suite - GOMo: <https://meme-suite.org/meme/tools/gomo>

❗ Image Not Showing

Possible Reasons

- The image was uploaded to a note which you don't have access to
- The note which the image was originally uploaded to has been deleted

[Learn More →](#)