# Calliphoridae Phylogeny

This notebook is for the description of the steps taken for the inference of a phylogeny for calliphorid species.

SPECIES (N = 15)

| Species | Genome | Transcriptome | Pedro | Vanessa |
|---------|--------|---------------|-------|---------|
| Agra | X | | x | |
| Bpan | X | | x | |
| Chom | X | | x | x |
| Cvom | X | | x | |
| Lcup | X | | x | x |
| Lser | X | | x | x |
| Pazu | X | | x | |
| Cput | | X | x | x |
| Lexi | | X | x | |
| Cmeg | | X | x | x |
| Calb | | X | | x |
| Clop | | X | | x |
| Loch | | X | | x |
| Cbez | | X | | x |
| Cmac | | X | | x |

*Cmac was added for Vanessa's study

# 0 - Promagrams

| Program | Version |
|---------|---------|
| TransDecoder | 5.7.0 |
| | |

| | |
|---|---|
| cd-hit | 4.8.1 |
| BUSCO | 5.4.7 |
| MAFFT | 7.505 |
| Trimal | 1.4.1 |
| IQTree | 2.2.2.6 |
| FigTree | 1.4.4 |

# 1 - Getting things that were ready before

We used the proteome from the species with genomes, and the transcriptomes from the othe ones. The BUSCO search for the proteomes was ready (Carol did it). We ran a new one for the proteins predicted from the transcriptomes.

```
1   # From the genomes
2   mkdir home/cunha/Phylogeny/0-Original_sequences
3
4   scp -r -P 2205 pedro@143.107.244.181:/home/blowflies/Gene_families/02-CDS_and
5
6   mkdir home/cunha/Phylogeny/1-BUSCO
7   scp -r -P 2205 pedro@143.107.244.181:/home/blowflies/Gene_families/04-BUSCO/1
8
9   for i in *busco; do mv $i ${i%%_protein_longest_isof.fasta.busco}_busco; done
10
11  # From the transcriptomes
12  cd /home/cunha/Phylogeny/0-Original_sequences
13
14  cp /home/cunha/Genotype_Phenotype/3-CDS/*.pep .
15
16  for i in *pep; do mv $i ${i%%.fasta.transdecoder.pep}_protein.fasta; done
17
18  for i in *trinity*; do sed -i -r 's/ .*//' $i; done
```

**OBS.: TransDecoder (Lexi)**

We did not have a "proteome" for Lexi, so we ran TransDecoder (the same way that was done for the other transcriptomes)

```
1   mkdir /home/cunha/Phylogeny/Lexi_transdecoder
2
3   cp /home/cunha/03-RNA/02-Trinity/lexi_trinity.fasta .
4
5   ~/Programs/TransDecoder-TransDecoder-v5.7.0/TransDecoder.LongOrfs -t lexi_tr
6
7   cd lexi_trinity.fasta.transdecoder_dir/
8
9   cp longest_orfs.pep /home/cunha/Phylogeny/0-Original_sequences/lexi_trinity_
```

Removing unnecessary parts of the id

```
1   for i in *; do sed -i -r 's/ .*//' $i; done
```

# 2 - cdhit (Transcriptomes)

We use this to remove the redundancy from the peptide files that came from the transcriptomes.

```
1   for i in *trinity*; do cd-hit -i $i -o $i.cdhit -c 0.97 -n 5 -T 5; done
2
3   # Renaming
4   for i in *trinity*fasta; do mv $i $i.NO-CD-HIT; done
5
6   for i in *cdhit; do mv $i ${i%%.cdhit}; done
```

# 3 - BUSCO (Transcriptomes)

```
1   cd home/cunha/Phylogeny/0-Original_sequences
2
3   for i in *trinity*fasta; do docker run -u $(id -u) -v $(pwd):/busco_wd ezlabg
4
5   mv *busco ../1-BUSCO
```

# 4 - Selection of Orthologs

Before running the command line, it is necessary to rename every full_table file with the first 3 letters or numbers of the respective sample

> Step 1

This command will select every Complete gene in BUSCO within the full_table file and assemble a new table with them.

```
cd /home/cunha/Phylogeny/1-BUSCO

# create-complete-txt.sh
for d in *_busco/
do
    echo "$d"
    cd "$d"
    cd run*/
        for i in $(find . -name full*)
            do grep Complete $i >> ${i%%.tsv}_complete.txt
            done
    cd ~/Phylogeny/1-BUSCO
done
```

## Step 2

Add the first 4 letters or numbers of the samples in each line of the complete genes

```
for i in $(find . -name *_complete.txt); do awk '{print
substr(FILENAME,3,4),$0}' $i > ${i%%.txt}_named.txt; done
```

## Step 3.1 (TRANSCRIPTOMA)

Add the samples names in the trinity files

```
cd /home/cunha/Phylogeny/0-Original_sequences

# name-in-trinity.sh
BEGIN {FS="\>"}
{
    if ($2 ~ /TRINITY_/) {
            print ">" substr(FILENAME,1,4) "_" $2
    } else {
            print $0
    }
}
```

And run the following command line

```
for i in *trinity*fasta ; do awk -f name-in-trinity.sh $i >>
/home/cunha/Phylogeny/2-Named_oneline_sequences/${i%%.fasta}_named.fasta; done
```

## Step 3.2 (GENOMAS)

Add the samples names in the genome files:

```
1  ls *_gen_* | sed "s/_gen_protein.fasta//g" > list.txt
```

Create file name_in_genome.sh with this code:

```
while read file
    do
    sed "s/>/>$file/g" ${file}_gen_protein.fasta >> /home/cunha/Phylogeny/2-
Named_oneline_sequences/${file}_gen_protein_named.fasta
done < list.txt
```

Transforming sequences in one line files

```
1   cd /home/cunha/Phylogeny/2-Named_oneline_sequences
2
3   for file in *; do awk '/^>/ {printf("\n%s\n",$0);next; } { printf("%s",$0);}
4
5   for i in *oneline*; do sed -i -e 1d $i; done #remove first line
6
7   rm *named.fasta
```

Select in the Hmmer output table the busco genes with evalue 0 or less than e^-100

```
1    cd /home/cunha/Phylogeny/1-BUSCO
2
3    for i in *_busco; do echo $i; done > dir.txt
4
5    # hmmr.sh (evalue < e-100)
6    while read file
7    do
8        SPP=$(echo $file | sed -r 's/_.*//')
9
10       cd /home/cunha/Phylogeny/1-BUSCO/$file/run_diptera_odb10/hmmer_output/in
11
12       for i in *.out*; do awk -v SPP2=$SPP '(NR == 4 && ($7 ~ "[0-9].[0-9]e-[0-
13
14        cd /home/cunha/Phylogeny/1-BUSCO
15
16   done < dir.txt
17
18   mv *evalue* ../3-evalues/
```

Assembling the final fasta files

```
1   cd /home/cunha/Phylogeny
2
3   # find - finds the full table with complete buscos with the species names
4   for i in $(find . -name *_complete_named.txt); do awk '{print $1 ".*" $4 ".*'
5
6   #find the genes with the evalue we selected
7   for i in $(find . -name *_grepfile.txt); do grep -hof $i $(find . -name *eval
8
9   # The files above appeared in the busco folders for each species
10
11  # This writes the command line to assemble the fasta files
12  cd /home/cunha/Phylogeny/2-Named_oneline_sequences
13
14  for i in $(find .. -name *_final.txt); do awk '{print "grep -h -A1 " $1 ".*"
15
16  screen
17  sh assemble_fasta_files.sh
18
19  mv *fas ../4-Orthologs_fasta
```

## Step 8

Filtering the orthogroups

```
1   cd /home/cunha/Phylogeny/4-Orthologs_fasta
2
3   mkdir All_ortho
4   mv *fas All_ortho/
5
6   mkdir min_7
7   mkdir min_8
8   mkdir min_9
9   mkdir min_10
10  mkdir min_11
11
12  #for Cmac
13  mkdir min_12
14
15  # filter.sh - filter by the number of sequences/species
16  for i in All_ortho/*
17  do
18      n=$(grep -c ">" $i)
19
20      if [ $n -gt 11 ]
21      then
22          cp $i min_12
23      fi
24  done
```
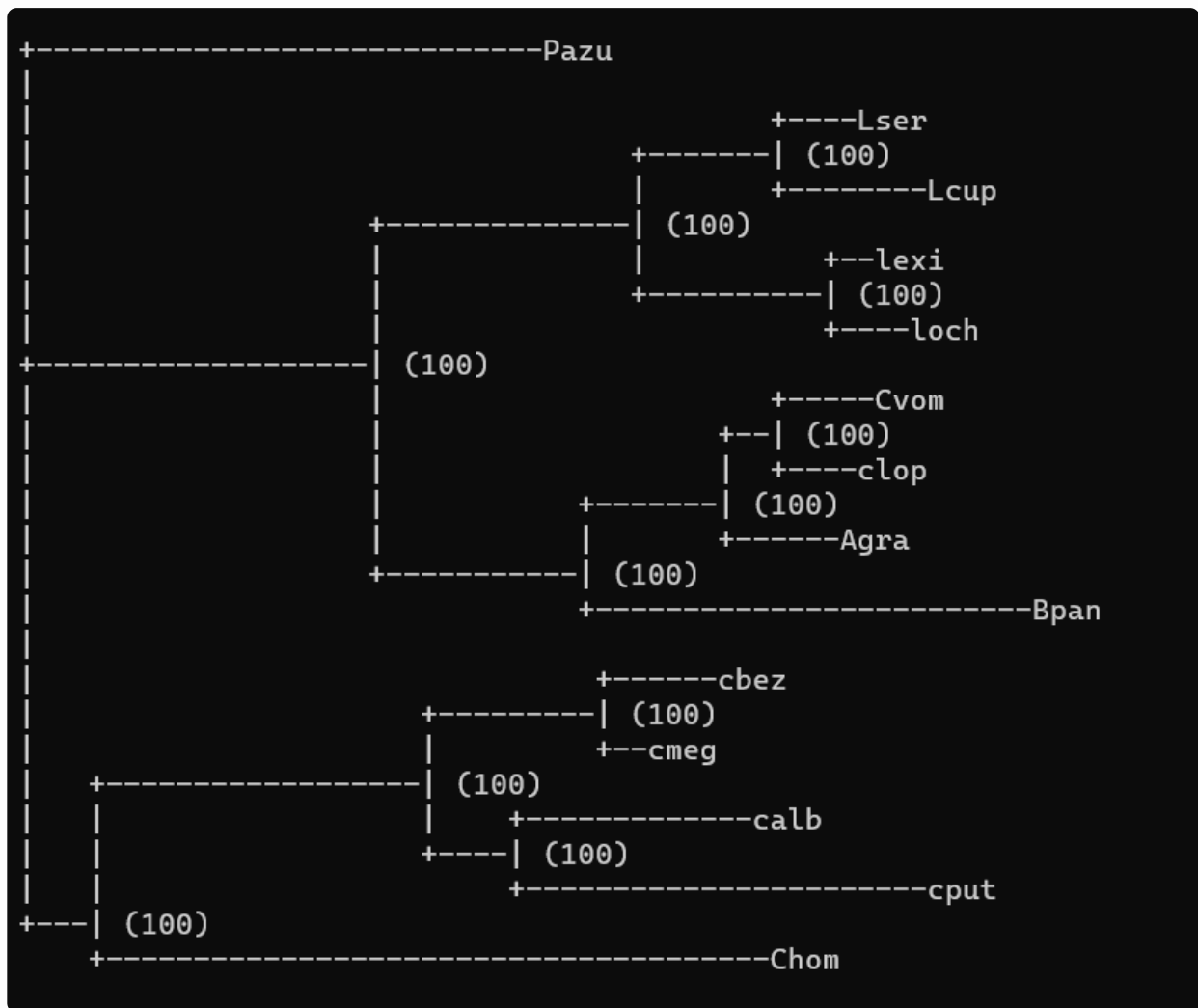
- min_12: XXXX (cov 75%)

We will infer the tree with the 75% cov.
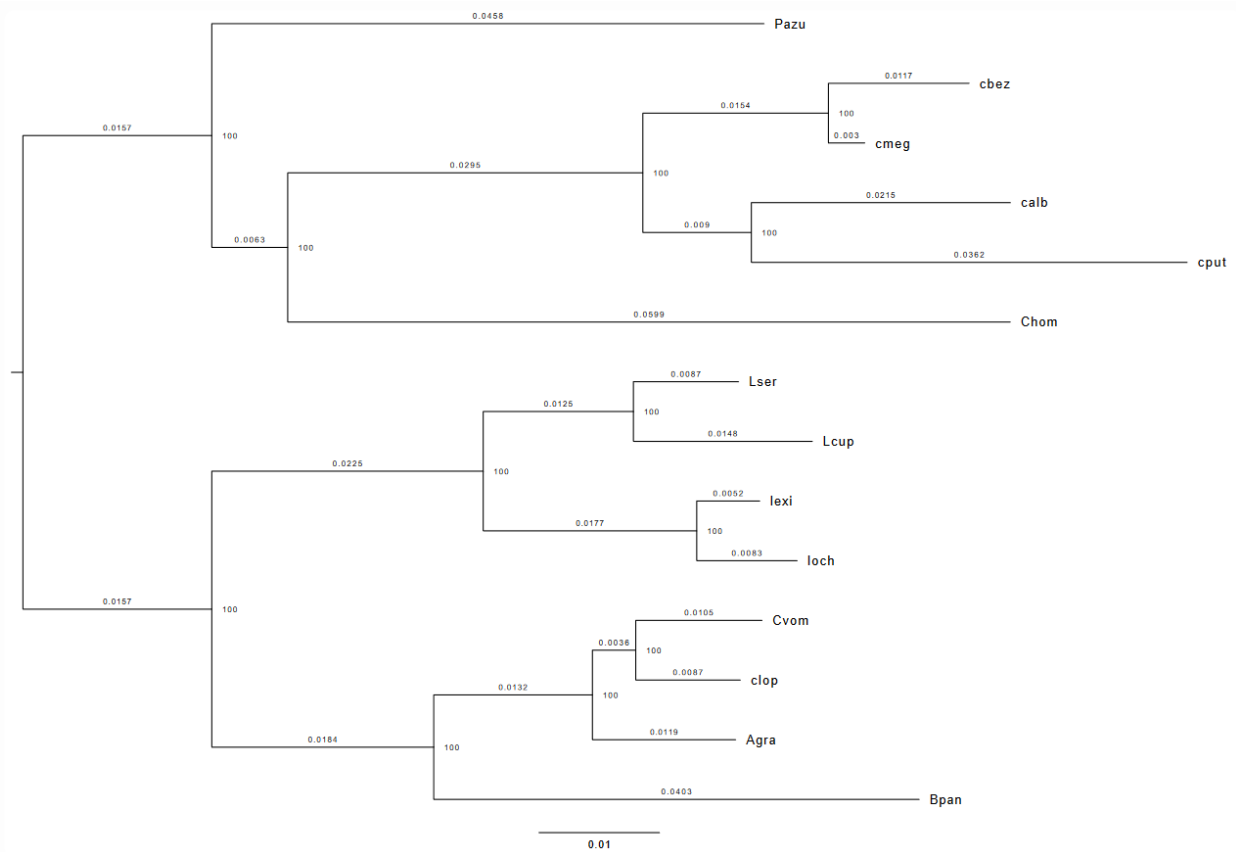
## 5 - Alignment

```
1   mkdir /home/cunha/Phylogeny/5-Alignments
2
3   cd /home/cunha/Phylogeny/4-Orthologs_fasta/min_11/
4
5   # alignment
6   for i in *fas; do mafft-ginsi --thread 25  --maxiterate 1000 --adjustdirecti
7
8   # cleaning the alignment
9   cd /home/cunha/Phylogeny/5-Alignments
10
11  for i in *fas; do trimal -in $i -out ${i%%.fas}_trim.fas -automated1; done
12
13  # fixing the sequence ids (leave just the species name)
14  for i in *trim*; do cat $i | sed 's/>Pazu.*/>Pazu/g' | sed 's/>Lser.*/>Lser/
15
16  mkdir 1-Mafft 2-Trimmed 3-Renamed
17  mv *renamed.fas 3-Renamed/
18  mv *trim.fas 2-Trimmed/
19  mv *aln.fas 1-Mafft/
```

## 6 - Tree inference

```
1   cd /home/cunha/Phylogeny/5-Alignments
2
3   iqtree2 -T 30 -s /home/cunha/Phylogeny/5-Alignments/3-Renamed/ -m MFP -bb 100
4
5   # the output files were named like the directory, so it's better to rename th
6   for i in 3-Renamed.*; do mv $i Calliphoridae_tree.${i##3-Renamed.}; done
7
8   mv Calliphoridae* /home/cunha/Phylogeny/6-IQ-TREE
```

```
+----------------------------------Pazu
|
|                                        +----Lser
|                               +-------| (100)
|                               |        +--------Lcup
|                 +------------| (100)
|                 |             |               +--lexi
|                 |             |        +---------| (100)
|                 |             |        |        +----loch
+------------------| (100)
|                 |                       +-----Cvom
|                 |                    +--| (100)
|                 |                    |  +----clop
|                 |             +-------| (100)
|                 |             |        +------Agra
|                 +----------| (100)
|                             +-----------------------Bpan
|
|                                        +------cbez
|                             +---------| (100)
|                             |          +--cmeg
|                +-----------------| (100)
|                |            |         +------------calb
|                |            +----| (100)
|                |                  +--------------------cput
+---| (100)
    +----------------------------------Chom
```

The tree matches the Dimensions phylogeny and a recently published one as well (Yan et al., 2021), with Chrysominae as a sister clade to Calliphorinae + Luciliinae. However, IQTREE infers an unrooted tree, so we rooted it in a way that the Pazu branch forms a clade with the other Chrysominae, using FigTree

Trees are here:

```
1   # Unrooted
2   /home/cunha/Phylogeny/6-IQ-TREE/Calliphoridae_tree.treefile
3   # Rooted
4   /home/cunha/Phylogeny/6-IQ-TREE/Calliphoridae_tree_root.tre
```