

# Calliphoridae coding-sequence evolution

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

## Table of contents

- Calliphoridae coding-sequence evolution
  - Table of contents
  - Programs (and prerequisites)
  - Working directory
  - 1. Genomic CDS
  - 2. Isoform removal
  - 3. Extraction of orthogroups
    - 3.1. Orthology assessment
    - 3.2 Paralog removal
  - 4. Filtering and alignment
    - 4.1. Filtering and fixing
    - 4.2. TranslatorX
  - 5. Test of evolutionary models
    - 5.1. Prepare data for CodeML
    - 5.2 Running CodeML
    - 5.3 Filtering
  - 6. GO analysis
    - 6.1 Annotation of GOs
    - 6.2 Comparison between the omegas associated to GOs and the average omega of the genome
    - 6.3 Enrichment analysis
  - 7. Extra analysis - Isoform removal by AGAT improves results
  - 8. Scripts
    - 8.1 "ete3\_wrapper\_.py": runs CodeML (Topic 5.1)
    - 8.2 "ete3\_parser.py": organizes important results from CodeML into tables (Topic 5.3)
    - 8.3 "codeml\_table\_filter.R": corrects p-values using FDR and compile genes for which the bestfit model was M0, bfree or bfree (Topic 5.3)
    - 8.4 "bind-tables\_go-omegas.R": binds tables of GOs and the omegas of their respective associated gene (Topic 6.2) (by ChatGPT)
    - 8.5 "MW\_test\_multipleX.R": performs Mann-Whitney U test with GOs that have at least X genes associated to them (Topic 6.2) (by ChatGPT)
    - 8.6 "MW\_test\_table.R": creates a table that compiles info necessary to plot a multiple boxplot graph (by ChatGPT)
    - 8.7 "GO\_terms.py": inputs a GO id list and outputs a table with ids and term names (by ChatGPT)
    - 8.8 "GO\_boxplot.R": script to create figures like the one from Scott et al. (2020), for the rate of evolution of the genes associated to specific GOs
    - 8.9 "get\_ids\_from\_cluster\_list.py"

## Programs (and prerequisites)

Program	Version
AGAT	v1.2.1

seqkit	v2.2.0
Orthofinder	v2.5.4
Phylopypruner	v1.2.4
Mafft	v7.505
TranslatorX	v1.1
ETE3	v3.1.3
GOATOOLS	v1.4.12

## Working directory

/home/blowflies/CDS\_evolution

### 1. Genomic CDS

For this analysis, we are going to use the genomic coding sequences (CDS) of the Calliphoridae species for which genomic data is available (Table 1).

Species	Accession number / Repository	# CDS
<i>Aldrichina grahami</i>	<a href="http://dx.doi.org/10.5524/100673">http://dx.doi.org/10.5524/100673</a>	12823
<i>Bellardia pandia</i> .	GCA_916048285.2	22397
<i>Cochliomyia hominivorax</i>		19678
<i>Chrysomya megacephala</i>	Data not published yet (Torres Lab)	27229
<i>Chrysomya putoria</i>	Data not published yet (Torres Lab)	20633
<i>Chrysomya vomitoria</i>	GCA_942486065.1	12917
<i>Lucilia cuprina</i>		13927
<i>Lucilia eximia</i>	Data not published yet (Torres Lab)	27610
<i>Lucilia sericata</i>		14704
<i>Protocalliphora azurea</i>		11936

### 2. Isoform removal

As orthologs of all species will be analysed, in order to reduce noise in the clusterization step of orthogroups, we will keep only the longest isoform per gene. This was carried out by Carolina when processing data for her Master's. Here is the link to her protocol: [https://hackmd.io/@canettieri/r1IK\\_aBtT](https://hackmd.io/@canettieri/r1IK_aBtT). The isoform removal is explained in the 2.2 topic.

Briefly, the removal of short isoforms was done using scripts from AGAT tool. It filters for longest isoforms based on information contained in annotation files as GFF/GTF. As it is explained in its user guide:

"The script aims to filter isoforms when present. For a locus: - when all isoforms have CDS we keep the one with the longest CDS. - when some isoforms have CDS some others not, we keep the one with the longest CDS. - when none of the isoforms have CDS, we keep the one with the longest concatenated exons."

### 3. Extraction of orthogroups

#### 3.1. Orthology assessment

```

1  mkdir /home/blowflies/CDS_evolution/1-Orthology/
2  cd /home/blowflies/CDS_evolution/1-Orthology/
3
4  mkdir WOI/
5  mkdir WOI/cds_WOI_data/
6
7  # Create a list with names of the species under analysis, i.e. Agra, Bpan, Cr
8  nano species-list
9
10 # Copy fasta files for cds with only the longest isoform
11 while read line; do rsync -c /home/blowflies/Gene_families/00-Genomic_data/0
12
13 # Copy Calliphoridae phylogenetic tree predicted
14 ## md5 was checked
15 ## Modify path to directory where the phylogeny was actually predicted (check
16 scp -P 2205 /home/martins/CODEML/Cactus_tree.tre carol@143.107.244.181:/home
17
18 # Edit fasta files' names to match to the species' names in the tree file, i
19
20 cd WOI/cds_WOI_data/
21 for file in *fasta; do mv $file ${file%%_cds_*}.fasta; done
22 for file in *fasta
23 do
24     NAME=${file%%.fasta}
25     sed -i "s/>/>${NAME}@/g" $file
26 done
27
28
29 # Run Orthofinder to predict orthology relationships
30 ## tmux a -t 3
31 /home/pedro/Programs/OrthoFinder/orthofinder -d -f WOI/cds_WOI_data/ -s Cactu

```

## 3.2 Paralog removal

### 3.2.1 Selection of orthogroups

We selected orthogroups containing at least 4 sequences, which is "the minimum number required for tree inference with most tree inference programs"

(<https://github.com/davideemms/OrthoFinder#gene-trees-directory>). Both sequences alignment and respective tree of sequences are the input for PhyloPyPruner.

```

1  cd /home/blowflies/CDS_evolution/1-Orthology/WOI
2  mkdir Phylopypruner/
3  mkdir Phylopypruner/Input/
4  mkdir Phylopypruner/Sequence_alignment/
5
6  cd Phylopypruner/
7
8  ## Copy gene trees and sequences of orthogroups
9  cp /home/blowflies/CDS_evolution/1-Orthology/WOI/Orthofinder/Results_Mar23/G
10 ls Input/ | wc -l # 14594
11
12 cd Input/
13 ls * | sed "s/_tree.txt//g" > ../orthogroups-names
14 cd ..
15 while read line; do cp /home/blowflies/CDS_evolution/1-Orthology/WOI/Orthofi

```

### 3.2.2 Alignment of sequences from orthogroups

```

1  # We did this step on Darwin
2  mkdir /home/martins/Mafft
3
4  # For WOI orthogroups
5  ## Align sequences from orthogroup
6  mkdir /home/martins/Mafft/WOI
7  cd /home/martins/Mafft/WOI
8
9  rsync -c -e "ssh -p 2205" carol@143.107.244.181:/home/blowflies/CDS_evolution
10
11 ## We ran Mafft through slurm. Below are the options used and written in WOI
12 #
13 ### =====
14 #!/usr/bin/bash
15
16 #SBATCH --job-name WOI ## nome que aparecerá na fila
17 #SBATCH --output=WOI.out ## nome do arquivo de saída; o %j é igual a jobID; o %J é igual a jobID
18 #SBATCH --ntasks=1 ## número de tarefas (análises) a serem executadas
19 #SBATCH --cpus-per-task=20 ## o número de threads alocados para cada tarefa
20 #SBATCH --partition=long ## as partições a serem executadas (separadas por vírgula)
21 #SBATCH --error=WOI.err
22
23 cd /home/martins/Mafft/WOI/
24
25 for file in *fa; do mafft --maxiterate 1000 --localpair --thread 20 $file > $file.out; done
26 ### =====
27 #
28 ## Run code in slurm
29 sbatch WOI.slurm

```

### 3.2.3 Running PhyloPyPruner

```

1  # Transfer aligned sequences from Darwin server to Rosalind server's input dir
2  ## The next two command lines were run on Darwin
3  cd /home/martins/Mafft/WOI
4  rsync -c -e "ssh -p 2205" *ali carol@143.107.244.181:/home/blowflies/CDS_evolution
5
6  # Rename tree and sequences files to match each other except by their extensions
7  cd /home/blowflies/CDS_evolution/1-Orthology/WOI/Phylopypruner/Input/
8
9  for file in *ali; do mv $file ${file%%.fa.ali}.fasta; done
10 for file in *txt; do mv $file ${file%%_tree.txt}.tree; done
11
12 # Rename taxa in trees to match headers in fasta
13 while read line
14 do
15     for file in *tree
16     do
17         sed -i "s/${line}_${line}@/${line}@/g" $file
18     done
19 done < /home/blowflies/CDS_evolution/1-Orthology/species-list
20
21 # Run phylopypruner
22 cd /home/blowflies/CDS_evolution/1-Orthology/WOI/Phylopypruner
23
24 /home/vanessa/anaconda3/bin/phylopypruner --dir /home/blowflies/CDS_evolution

```

### 3.2.4 Results

Alignment statistics:		
Description	Input	Output
No. of alignments	14594	11322
No. of sequences	146330	98395
No. of OTUs	10	10
Avg no. of sequences / alignment	10	8
Avg no. of OTUs / alignment	7	8
Avg sequence length (ungapped)	1664	1790
Shortest sequence (ungapped)	54	54
Longest sequence (ungapped)	72516	72516
% missing data	32.60	23.80
Concatenated alignment length	30646102	24314468

- Number of 1:1 orthogroups from Orthofinder = ?
- Number of orthogroups with at least 4 species from Orthofinder (our input): 14594
- Number of orthogroups without paralogs from Phylopypruner: 11322
- Number of input sequences = 146,330
- Number of filtered sequences = 98,395
- Percentage of filtered sequences in relation to input sequences = 67,24%



```

1 | mkdir /home/martins/CodeML
2 | mkdir /home/martins/CodeML/Chom /home/martins/CodeML/Pazu
3 | cd /home/martins/CodeML
4 |
5 | # Tree
6 | ## Modify path to directory where the phylogeny was actually predicted (check)
7 | rsync -c CODEML/Cactus_tree.tre .
8 |
9 | # Aligned sequences
10 | rsync -c -e "ssh -p 2205" carol@143.107.244.181:/home/blowflies/CDS_evolution
11 | rsync -c -e "ssh -p 2205" carol@143.107.244.181:/home/blowflies/CDS_evolution
12 |
13 | # Scripts
14 | nano CodeML/Chom/ete3_wrapper_Chom.py
15 | nano CodeML/Pazu/ete3_wrapper_Pazu.py

```

See scripts in Topic 9.

## 5.2 Running CodeML

We ran CodeML using the following commands in slurm.

For Chom:

```

1 | #!/usr/bin/bash
2 | #SBATCH --job-name Chom_codeml ## nome que aparecerá na fila
3 | #SBATCH --output Chom_codeml.out ## nome do arquivo de saída; o %j é igual a
4 | #SBATCH --ntasks=1 ## número de tarefas (análises) a serem executadas
5 | #SBATCH --cpus-per-task=30 ## o número de threads alocados para cada tarefa
6 | #SBATCH --partition=long ## as partições a serem executadas (separadas por v:
7 | #SBATCH --error=Chom_codeml.err
8 |
9 | cd /home/martins/CodeML/Chom
10 | python3 ete3_wrapper_Chom.py -p /home/martins/CodeML/Chom/Aln -e .fasta -t /l

```

For Pazu:

```

1 | #!/usr/bin/bash
2 | #SBATCH --job-name Pazu_codeml ## nome que aparecerá na fila
3 | #SBATCH --output Pazu_codeml.out ## nome do arquivo de saída; o %j é igual a
4 | #SBATCH --ntasks=1 ## número de tarefas (análises) a serem executadas
5 | #SBATCH --cpus-per-task=30 ## o número de threads alocados para cada tarefa
6 | #SBATCH --partition=long ## as partições a serem executadas (separadas por v:
7 | #SBATCH --error=Pazu_codeml.err
8 |
9 | cd /home/martins/CodeML/Pazu
10 | python3 ete3_wrapper_Pazu.py -p /home/martins/CodeML/Pazu/Aln -e .fasta -t /l

```

sbatch Chom\_codeml.slurm

sbatch Pazu\_codeml.slurm

## 5.3 Filtering

We organized the results into tables with the important information using another script.

```

1 | cd /home/martins/CodeML
2 |
3 | rsync -c -e "ssh -p 2205" carol@143.107.244.181:/home/pedro/Calliphoridae_pe
4 |
5 | python3 ete3_parser.py -p Chom -e .log -s Chom
6 | python3 ete3_parser.py -p Pazu -e .log -s Pazu

```

See script in Topic 9.

After that, we corrected p-values using FDR in R and selected genes for which the bestfit model was MO, b\_neut and b\_free.

```

1 | mv Chom/Chom_branch_raw_table.txt .
2 | mv Pazu/Pazu_branch_raw_table.txt .
3 |
4 | Rscript codeml_table_filter.R
5 |
6 | rsync -c -e "ssh -p 2205" -a /home/martins/CodeML/ carol@143.107.244.181:/ho

```

See script in Topic 9.

We have 6 tables at the end (3 for each species):

Table	Content
[Chom/Pazu]_bestfit_M0.txt	Genes that did not reject the first null hypothesis (M0)
[Chom/Pazu]_bestfit_neut.txt	Genes that rejected the first null hypothesis (M0), but did not reject the second one (b_neut)
[Chom/Pazu]_bestfit_free.txt	Genes that rejected both null hypotheses and favoured the alternative one (free)

The following genes do not have significant omega because their aligned sequence (input for CodeML) were empty:

- Chom: "OG0000190\_pruned\_Chom", "OG0003185\_pruned\_1\_Chom", "OG0003142\_pruned\_1\_Chom"
- Pazu: "OG0000190\_pruned\_Pazu", "OG0001498\_pruned\_1\_Pazu"

Thus, at the end, we have 9459 and 8924 genes with significant omegas for Chom and Pazu, respectively.

## 6. GO analysis

```
1 | mkdir /home/blowflies/CDS_evolution/4-GO_Analysis
2 | cd /home/blowflies/CDS_evolution/4-GO_Analysis
```

### 6.1 Annotation of GOs

We annotated only the genes from Chom and Pazu. We did this step on Darwin server using EnTAP software. Note that the headers are named as `$SPECIES@GENE_NAME`

```
1 | mkdir /home/martins/CDS_EnTAP_Pedro
2 | cd /home/martins/CDS_EnTAP_Pedro
3
4 | mkdir Chom Pazu
5
6 | ## Copy files
7 | rsync -c -e "ssh -p 2205" carol@143.107.244.181:/home/blowflies/CDS_evolution/4-GO_Analysis/
8 | rsync -c -e "ssh -p 2205" carol@143.107.244.181:/home/blowflies/CDS_evolution/4-GO_Analysis/
9
10 | # Run EnTAP with cds in queue using slurm
11
12 | cd Chom # slurm file: Chom.slurm
13 | EnTAP --runN -i /home/martins/CDS_EnTAP_Pedro/Transcripts/Chom.fasta -d /home/martins/CDS_EnTAP_Pedro/Transcripts/Chom.fasta
14
15 | cd Pazu # slurm file: Pazu.slurm
16 | EnTAP --runN -i /home/martins/CDS_EnTAP_Pedro/Transcripts/Pazu.fasta -d /home/martins/CDS_EnTAP_Pedro/Transcripts/Pazu.fasta
17
18 | # Transfer files to rosalind
19
20 | # In Rosalind previously
21 | ## mkdir /home/blowflies/CDS_evolution/4-GO_Analysis/EnTAP
22 | rsync -c -e "ssh -p 2205" /home/martins/CDS_EnTAP_Pedro/Chom/entap_outfiles/ /home/blowflies/CDS_evolution/4-GO_Analysis/EnTAP/
23 | rsync -c -e "ssh -p 2205" /home/martins/CDS_EnTAP_Pedro/Pazu/entap_outfiles/ /home/blowflies/CDS_evolution/4-GO_Analysis/EnTAP/
```

### 6.2 Comparison between the omegas associated to GOs and the average omega of the genome

The analysis onwards were run on Rosalind server

```
1 | mkdir /home/blowflies/CDS_evolution/4-GO_Analysis/2-Distribution_omega
2 | cd /home/blowflies/CDS_evolution/4-GO_Analysis/2-Distribution_omega
3
4 | # Write Chom and Pazu in a list
5 | nano species-list.txt
```

Create a table with genes and their omegas estimated by their bestfit model

```

1 while read sp
2 do
3     sed '1d' /home/bowflies/CDS_evolution/3-CodeML/${sp}_bestfit_M0.txt | cut -f1 > ${sp}_bestfit_M0.txt
4     sed '1d' /home/bowflies/CDS_evolution/3-CodeML/${sp}_bestfit_free.txt | cut -f1 > ${sp}_bestfit_free.txt
5     sed '1d' /home/bowflies/CDS_evolution/3-CodeML/${sp}_bestfit_neut.txt | cut -f1 > ${sp}_bestfit_neut.txt
6     cat ${sp}*txt | sed '1i\Gene\tOmega\tBestfit_Model' > ${sp}_omegas_bestfit.txt
7     rm ${sp}_bestfit*txt
8 done < species-list.txt
9
10 # Check
11 wc -l Chom* # 9460 (header included)
12 wc -l Pazu* # 8925 (header included)

```

Filter genes with omega from EnTAP table

```

1 # Get ID of orthogroups that were tested in CodeML
2 while read sp
3 do
4     sed '1d' ${sp}_omegas_bestfit.txt | cut -f1 | sed "s/_${sp}/g" > ${sp}_orthogroups.txt
5 done < species-list.txt
6
7 # Get IDs of genes that were tested in CodeML
8 cd /home/bowflies/CDS_evolution/4-GO_Analysis/1-EnTAP/All_orthogroups
9
10 while read sp
11 do
12     while read id
13     do
14         cat ${id}.fasta | grep "${sp}" | sed "s/>${sp}@/g" >> ../2-Distribution_omega/${sp}.tmp
15     done < /home/bowflies/CDS_evolution/4-GO_Analysis/2-Distribution_omega/species-list.txt
16 done < /home/bowflies/CDS_evolution/4-GO_Analysis/2-Distribution_omega/species-list.txt
17
18 # Check
19 wc -l /home/bowflies/CDS_evolution/4-GO_Analysis/2-Distribution_omega/*genes.tmp
20
21 # Filter EnTAP table for those genes
22 cd /home/bowflies/CDS_evolution/4-GO_Analysis/EnTAP
23
24 while read sp
25 do
26     grep -f /home/bowflies/CDS_evolution/4-GO_Analysis/2-Distribution_omega/${sp}.tmp EnTAP.tsv > ${sp}.tmp
27     cut -f1 ${sp}.tmp | grep -v -f /home/bowflies/CDS_evolution/4-GO_Analysis/2-Distribution_omega/species-list.txt > ${sp}.del
28     grep -v -f ${sp}.del EnTAP.tsv | sort > /home/bowflies/CDS_evolution/4-GO_Analysis/2-Distribution_omega/species-list.txt
29     rm ${sp}.tmp ${sp}.del
30 done < /home/bowflies/CDS_evolution/4-GO_Analysis/2-Distribution_omega/species-list.txt
31
32 # Check
33 wc -l /home/bowflies/CDS_evolution/4-GO_Analysis/2-Distribution_omega/*genes.tmp
34
35 # Create tables with infos
36 cd /home/bowflies/CDS_evolution/4-GO_Analysis/2-Distribution_omega
37
38 while read sp
39 do
40     paste ${sp}_omegas_genes.id ${sp}_omegas_bestfit.txt | sort | paste - ${sp}_orthogroups.txt > ${sp}_all-infos.tsv
41     cut -f1,2,3,4,17,38,39,40 ${sp}_all-infos.tsv > ${sp}_short-infos.tsv
42 done < species-list.txt

```

Create a table with GOs and the omega of the genes associated to them



```

1  ## Create a list with GO categories (Write BP, CC and MF)
2  nano go-categories.txt
3
4  ## Get the GO terms present in the dataset per category
5  while read sp
6  do
7      while read go
8      do
9          if [[ $go == "BP" ]]
10         then
11             sed 'ld' ${sp}_short-infos.tsv | cut -f6 | sed -e $'s/,/\n/g'
12         elif [[ $go == "CC" ]]
13         then
14             sed 'ld' ${sp}_short-infos.tsv | cut -f7 | sed -e $'s/,/\n/g'
15         elif [[ $go == "MF" ]]
16         then
17             sed 'ld' ${sp}_short-infos.tsv | cut -f8 | sed -e $'s/,/\n/g'
18         fi
19     done < go-categories.txt
20 done < species-list.txt
21
22 ## Get the number of genes that each GO is associated to
23 N=10; i=0; while read go; do ((i=i%N)); ((i+=1)) && wait; (while read id; do
24
25 N=10; i=0; while read go; do ((i=i%N)); ((i+=1)) && wait; (while read id; do
26
27 ## Create a table with GOs and the omega of the genes associated to them
28 N=10; i=0; while read go; do ((i=i%N)); ((i+=1)) && wait; (mkdir Chom_${go}
29
30 N=10; i=0; while read go; do ((i=i%N)); ((i+=1)) && wait; (mkdir Pazu_${go}
31
32 # Bind tables into one
33 N=5; i=0; while read sp; do ((i=i%N)); ((i+=1)) && wait; (while read go; do
34
35 # Check
36 for file in *amount; do echo $file; sort -k2n $file | tail -n1; done
37 # Chom_BP.amount
38 # GO:0008150      7506
39 # Chom_CC.amount
40 # GO:0005575      6254
41 # Chom_MF.amount
42 # GO:0003674      7344
43 # Pazu_BP.amount
44 # GO:0008150      7169
45 # Pazu_CC.amount
46 # GO:0005575      5973
47 # Pazu_MF.amount
48 # GO:0003674      6993
49 wc -l *go-omegas.tsv # headers included
50 # 7507 Chom_BP_go-omegas.tsv
51 # 6255 Chom_CC_go-omegas.tsv
52 # 7345 Chom_MF_go-omegas.tsv
53 # 7170 Pazu_BP_go-omegas.tsv
54 # 5974 Pazu_CC_go-omegas.tsv
55 # 6994 Pazu_MF_go-omegas.tsv

```

Perform Mann-Whitney U test for GOs that have at least X omegas of genes associated to them and, among GOs whose median is significantly different from genomes's median, it select the top 15 GOs that have the highest and the lowest medians compared to the genome's median (top 30 GOs at the end).

The user can input multiple values of X to see how many significant GOs they get at the end

A graph is created to visualize the amount of significant GOs versus the X values specified (

num\_significant\_plot.png) The output table helps the user to visualize the top 30 GOs whose median are the highest and the lowest\_top15\_highest/lowest\_headers.tsv)

```

1  N=5; i=0; while read sp; do ((i=i%N)); ((i+=1)) && wait; (while read go; do
2
3  # Just a command to verify how many GOs have at least a X value of genes asso
4  awk -F "\t" '{ if ($2 > 500) print $1}' Chom_BP.amount | wc -l

```

When X = 30, a higher number of GOs were returned as the top 15 highest median. The results indicates that GOs associated to more genes, that is, the more generalist a GO are, the more close its median gets to genome's median. As we want to get more detailed information about the genome evolution of our species, we will proceed the analysis using X = 30.

Perform Mann-Whitney U test for GOs that have at least 30 omegas of genes associated to them, select the top 30 genes (see previously) and creates a table that compiles info necessary to plot a multiple boxplot graph (see Figure 5 of "Genomic analysis of a livestock pest (Scott et al., 2020)")

Output written in (\$\_top30\_GOs.tsv)

```
1 | N=5; i=0; ((i=i%N)); ((i+=1)) && wait; while read sp; do Rscript MW_test_t
```

Before going to the following steps, we had to do some manual changes to the tables (if not specified otherwise, everything was done on Excel):

- GO term tables (\*\_top30\_GOs.tsv):
  1. The GO ids were written with ".", so we changed it to ":", because that's how it is recognized by GO-related tools (ex: "GO.0050907" to "GO:0050907");
  2. We used a script made by ChatGPT to obtain the GO term names from the ids (see script 8.7), which returns a two-column table with the ids and the names;
  3. We manually inserted the name column into the original table;
  4. We removed the header.
- Whole genome table (\*\_short-infos.tsv):
  1. We removed everything but the column with all omegas.

These were the inputs used in the R script we used to generate the final figures (script 8.8)

See scripts in Topic 8

We did everything locally (because of RStudio) and uploaded everything back to Rosalind afterwards: `/home/blowflies/CDS_evolution/4-GO_Analysis/3-GO_figure`

## 6.3 Enrichment analysis

### Organization

We performed the enrichment analysis using GOATools. This software requires three inputs:

1. A file with a list of all genes in a given genome (population);
2. A file with all genes of interest (study - in our case, the genes whose best fit was the b\_free model);
3. An association file, which relates all GO-annotated genes to their GOs.

```

1  cd /home/blowflies/CDS_evolution/4-GO_Analysis/4-Enrichment
2
3  # Population files
4  /home/blowflies/CDS_evolution/1-Orthology/WOI/cds_WOI_data
5
6  grep ">" /home/blowflies/CDS_evolution/1-Orthology/WOI/cds_WOI_data/Pazu.fast
7  grep ">" /home/blowflies/CDS_evolution/1-Orthology/WOI/cds_WOI_data/Chom.fast
8
9  # Study files
10 ## These are a bit more challenging. Our codeml outputs contain the orthogroup
11
12 ## 1. Get the names of the orthogroups we need, and their omega values
13 cut -f1,5 /home/blowflies/CDS_evolution/3-CodeML/Chom_bestfit_free.txt | sed
14
15 cut -f1,5 /home/blowflies/CDS_evolution/3-CodeML/Pazu_bestfit_free.txt | sed
16
17 ## 2. Separate the orthogroups into the ones with above and below average (co
18
19 # Chom - mean = 0.073
20 while IFS=$'\t' read -r word value; do
21     if (( $(echo "$value >= 0.073" | bc -l) )); then
22         echo -e "$word" >> Chom_above.txt
23     elif (( $(echo "$value < 0.073" | bc -l) )); then
24         echo -e "$word" >> Chom_below.txt
25     fi
26 done < Chom_orthogroups.txt
27
28 # Pazu - mean = 0.082
29 while IFS=$'\t' read -r word value; do
30     if (( $(echo "$value >= 0.082" | bc -l) )); then
31         echo -e "$word" >> Pazu_above.txt
32     elif (( $(echo "$value < 0.082" | bc -l) )); then
33         echo -e "$word" >> Pazu_below.txt
34     fi
35 done < Pazu_orthogroups.txt
36
37 ## 3. Extract the names of the genes of each species from the orthogroups (tl
38 cp /home/pedro/Scripts/get_ids_from_cluster_list.py .
39
40 python3 get_ids_from_cluster_list.py -i Chom_above.txt -s Chom -c /home/blow
41 mv Chom_table_from_clusters.txt Chom_table_from_clusters_above.txt
42 cut -f2 Chom_table_from_clusters_above.txt | sed 's/>Chom@//' > Chom_above_s
43
44 python3 get_ids_from_cluster_list.py -i Chom_below.txt -s Chom -c /home/blow
45 mv Chom_table_from_clusters.txt Chom_table_from_clusters_below.txt
46 cut -f2 Chom_table_from_clusters_below.txt | sed 's/>Chom@//' > Chom_below_s
47
48 python3 get_ids_from_cluster_list.py -i Pazu_above.txt -s Pazu -c /home/blow
49 mv Pazu_table_from_clusters.txt Pazu_table_from_clusters_above.txt
50 cut -f2 Pazu_table_from_clusters_above.txt | sed 's/>Pazu@//' > Pazu_above_s
51
52 python3 get_ids_from_cluster_list.py -i Pazu_below.txt -s Pazu -c /home/blow
53 mv Pazu_table_from_clusters.txt Pazu_table_from_clusters_below.txt
54 cut -f2 Pazu_table_from_clusters_below.txt | sed 's/>Pazu@//' > Pazu_below_s
55
56 # Association files (preliminary - We need to do some more fixing on Excel!)
57 cut -f1,34,35,36 /home/blowflies/CDS_evolution/4-GO_Analysis/EnTAP/Chom_entap
58
59 cut -f1,34,35,36 /home/blowflies/CDS_evolution/4-GO_Analysis/EnTAP/Pazu_entap
60
61 # Changes in Excel:
62 # 1. Substitute "NA" for nothing
63 # 2. Concatenate the three GO columns in another one
64 # 3. Substitute ",GO" for ";GO"
65 # 4. Substitute "," for nothing (remove the final comma)
66 # 5. Reorder the table (Z -> A), so it's easier to remove the genes with no
67 # 6. For some reason, when we open the files in a txt reader, quotes (") app
68 # 7. Rename the files: Chom_id2GO.txt and Pazu_id2GO.txt (these are in the E

```

Just a little organization to make our lives easier afterwards

```

1  mkdir /home/blowflies/CDS_evolution/4-GO_Analysis/Enrichment/1-Initial_files
2  mv *txt 1-Initial_files
3
4  # Now, let's put the non-necessary files elsewhere, and leave only the neces
5  cd 1-Initial_files
6  mkdir TMP
7  mv mv Chom_above.txt Chom_below.txt *table* *orthogroups* *tmp* TMP/
8
9  # Done! Now we have one population file, two study files and one association

```

## Analysis

We started by creating all necessary files on Rosalind, and then moved on to analysing them on Darwin. So we had to install the software and move the files.

```
1 # On Darwin
2 mkdir /home/martins/GO_enrichment
3 cd /home/martins/GO_enrichment
4
5 # Installation (July, 13 - 2024 -> Reference for the GO obo file)
6 wget http://current.geneontology.org/ontology/go-basic.obo
7 pip install goatools
8
9 # Bring the files
10 scp -P 2205 pedro@143.107.244.181:/home/blowflies/CDS_evolution/4-GO_Analysis
11
12 # Analysis
13 find_enrichment.py Chom_above_study.txt Chom_pop.txt Chom_id2GO.txt --obsolete
14 find_enrichment.py Chom_below_study.txt Chom_pop.txt Chom_id2GO.txt --obsolete
15 find_enrichment.py Pazu_above_study.txt Pazu_pop.txt Pazu_id2GO.txt --obsolete
16 find_enrichment.py Pazu_below_study.txt Pazu_pop.txt Pazu_id2GO.txt --obsolete
17
18 # Note: the go-basic.obo file is in the same directory as all the files, so :
```

Bringing the raw results back to Rosalind

```
1 scp -P 2205 *enriched* pedro@143.107.244.181:/home/blowflies/CDS_evolution/4-
```

### Filtering and figures

We downloaded the GOATOOLS outputs and hadled them using R (GO\_enrichment\_filter\_plots.R - see 8.10). Afterwards, we brought everything back to Rosalind.

```
1 scp -P 2205 * pedro@143.107.244.181:/home/blowflies/CDS_evolution/4-GO_Analy:
```

## 7. Extra analysis - Isoform removal by AGAT improves results

When we first ran this pipeline, we tested the analysis with two datasets: one containing only the longest isoform for each gene and another containing all isoforms, for all species. By doing this, we aimed to evaluate if AGAT's script to remove short isoforms would increase the number of orthogroups clustered by Phylopypruner in comparison to the number the latter can cluster with all-isoforms-dataset using its option '-mask longest'. The results of Phylopypruner for each dataset are shown below:

### For dataset containig only the longest isoform (AGAT's script used)

Alignment statistics:		
Description	Input	Output
No. of alignments	14630	11347
No. of sequences	145793	98661
No. of OTUs	10	10
Avg no. of sequences / alignment	9	8
Avg no. of OTUs / alignment	7	8
Avg sequence length (ungapped)	1658	1794
Shortest sequence (ungapped)	54	57
Longest sequence (ungapped)	72588	72588
% missing data	32.80	23.90
Concatenated alignment length	30767442	24385049

- Number of 1:1 orthogroups from Orthofinder = 2336
- Number of orthogroups with at least 4 species from Orthofinder (our input): 14630
- Number of orthogroups without paralogs from Phylopypruner: 11347
- Number of filtered sequences = 98,661
- Percentage of filtered sequences in relation to input sequences = 67.7%

#### For dataset containing all isoforms (Phylopypruner's masking option used)

Alignment statistics:		
Description	Input	Output
No. of alignments	15515	11234
No. of sequences	182957	95704
No. of OTUs	10	10
Avg no. of sequences / alignment	11	8
Avg no. of OTUs / alignment	7	8
Avg sequence length (ungapped)	1811	1745
Shortest sequence (ungapped)	54	57
Longest sequence (ungapped)	72588	70389
% missing data	34.80	25.50
Concatenated alignment length	32314878	23800455

- Number of 1:1 orthogroups from Orthofinder = 752
- Number of orthogroups with at least 4 species from Orthofinder (our input): 15515
- Number of orthogroups without paralogs from Phylopypruner: 11234
- Number of filtered sequences = 95,704
- Percentage of filtered sequences in relation to input sequences = 52,3%

Looking at the results, we can say that using AGAT's script provided 110 more orthogroups and almost 3,000 more sequences to the analysis, which is advantageous. In this sense, we can recommend using this tool instead of just Phylopypruner's masking option.

Because Vanessa's detected a problem with the annotation of Cmeg, Lexi and Cput, we had to run all the pipeline again. As we had already run this comparison, it would not make sense to rerun it. We decided to directly use AGAT.

In order to not confound other users, we deleted the files regarding the results shown.

However, we preserved their statistics in this section in order to convince new users to use AGAT.

## 8. Scripts

### 8.1 "ete3\_wrapper\_.py": runs CodeML (Topic 5.1)

The script below is for Pazu, but works the same for Chom if you substitute Pazu for Chom:

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  # ete3_wrapper.py
4
5  import sys
6  import os
7  import getopt
8  import glob
9  from ete3 import EvolTree
10 #-----
11 def main():
12     pwd, end, t = arg()
13
14     if pwd == None or end == None or t == None:
15         usage()
16
17     clusters = glob.glob(f"{pwd}/{end}")
18
19     for cluster in clusters:
20         # Get the name for the output directory
21         name = cluster.split("/")[-1].split(".")[0]
22
23         # What species are represented in the fasta file?
24         species = []
25
26         fasta = open(cluster, 'r', encoding='utf-8')
27
28         for line in fasta:
29             if ">" in line:
30                 sp = line.split(">")[1].strip()
31                 species += [sp]
32
33         fasta.close()
34
35
36         # Read the tree and prune it leaving just the species that have sequ
37         tree = EvolTree(t)
38         tree.prune(species)
39
40         tree_txt = tree.write()
41
42         # New tree
43         new_tree = open('temp_tree.nw', 'w')
44         new_tree.write(tree_txt)
45         new_tree.close()
46
47         final_tree = "temp_tree.nw"
48
49
50         # Run CODEML with ete3
51         cmd = ["ete3 evol -t", final_tree, "--alg ", cluster, "--mark Pazu", "\
52             \"--models b_free b_neut M0 --tests b_free,b_neut b_free,M\
53             f\"-o {name}_Pazu/ >{name}_Pazu.log 2>{name}_Pazu.err"]
54         print (" ".join(cmd))
55         os.system(" ".join(cmd))
56
57 #-----
58 def arg():
59     argv = sys.argv[1:]
60     try:
61         opts, args = getopt.getopt(argv, 'p:e:t:')
62
63     except:
64         usage()
65
66     arg_1, arg_2, arg_3 = None, None, None
67
68     for opt, arg in opts:
69         if opt in ['-p']:
70             arg_1 = arg
71         elif opt in ['-e']:
72             arg_2 = arg
73         elif opt in ['-t']:
74             arg_3 = arg
75
76     return arg_1, arg_2, arg_3
77 #-----
78 def usage():
79     sys.exit('Usage:\n'
80             '\t-p: fasta clusters directory pathway\n'
81             '\t-e: fasta clusters name ending\n'
82             '\t-t: tree file (newick format)')
83 #-----
84 if __name__ == '__main__':
85     main()

```

## 8.2 "ete3\_parser.py": organizes important results from CodeML into tables (Topic 5.3)

For reference of ete3\_parser.py:





```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3
4  '''
5  PROGRAM: ete_codeml_parser.py
6
7  DESCRIPTION: Extracts CodeML results from many ete3 result files using regex
8  '''
9  import re
10 import glob
11 import sys
12 import getopt
13 #-----
14 def main():
15     pwd, end, spp = arg()
16     if pwd == None or end == None or spp == None:
17         usage()
18
19     results = glob.glob(f'{pwd}/{end}')
20
21     species = spp.split(',')
22
23     for s in species:
24         # two tables per species
25         table_branch = open(f'{pwd}/{s}_branch_raw_table.txt','w')
26         table_branch.write('Gene\tAverage_Omega_M0\tOmega_foreground_neut\t'
27                             'Omega_background_neut\tOmega_foreground_free\t'
28                             'Omega_background_free\tpv_free_M0\tpv_free_neut\t'
29
30     for r in results:
31         if s in r: # just the results for that species
32             file = open(r,'r')
33             gene = r.split('/')[1].split('.')[0]
34
35             x = False # variable to check where we are in the file (beca
36                         # there two lines that are equal to each other, bu
37                         # need separate information -> this is for omega f
38                         # and background values)
39
40             pv_free_M0 = 'NA'
41             pv_free_neut = 'NA'
42             omega_fore_neut = 'NA'
43             omega_back_neut = 'NA'
44             omega_fore_free = 'NA'
45             omega_back_free = 'NA'
46             av_omega_M0 = 'NA'
47
48             for line in file:
49                 # pvalue for branch test free x M0
50                 pv_free_M0_search = re.search(r'M0.*? b_free.*?\| (.*)',
51                 if pv_free_M0_search:
52                     pv_free_M0 = pv_free_M0_search.group(1)
53                     pv_free_M0 = pv_free_M0.replace('*', '')
54                     continue
55
56                 # pvalue for branch test free x neut
57                 pv_free_neut_search = re.search(r'b_neut.*? b_free.*?\|
58                 if pv_free_neut_search:
59                     pv_free_neut = pv_free_neut_search.group(1)
60                     pv_free_neut = pv_free_neut.replace('*', '')
61                     continue
62
63                 # this line here separates the file
64                 if '- Model b_neut' in line:
65                     x = True
66                     continue
67
68                 # get omega values, but we differentiate the line by whe
69                 # x is True or False
70                 omega_fore_search = re.search(r'#1 => (.*)', line)
71                 if omega_fore_search and not x:
72                     omega_fore_free = omega_fore_search.group(1)
73                     omega_fore_free = omega_fore_free.replace(' ', '')
74                     continue
75                 elif omega_fore_search and x:
76                     omega_fore_neut = omega_fore_search.group(1)
77                     omega_fore_neut = omega_fore_neut.replace(' ', '')
78                     continue
79
80                 omega_back_search = re.search(r'background => (.*)',
81                 if omega_back_search and not x:
82                     omega_back_free = omega_back_search.group(1)
83                     omega_back_free = omega_back_free.replace(' ', '')
84                     continue
85                 elif omega_back_search and x:
86                     omega_back_neut = omega_back_search.group(1)
87                     omega_back_neut = omega_back_neut.replace(' ', '')
88                     continue

```

```

89
90         av_omega_search = re.search(r'Average omega for all tree
91         if av_omega_search:
92             av_omega_M0 = av_omega_search.group(1)
93             av_omega_M0 = av_omega_M0.replace(' ','')
94             continue
95
96         table_branch.write(f'{gene}\t{av_omega_M0}\t{omega_fore_neut
97                           f'{omega_back_neut}\t{omega_fore_free}\t'
98                           f'{omega_back_free}\t{pv_free_M0}\t'
99                           f'{pv_free_neut}\n')
100
101     table_branch.close()
102
103
104 #-----
105 def arg():
106     argv = sys.argv[1:]
107     try:
108         opts,args = getopt.getopt(argv, 'p:e:s:')
109     except:
110         usage()
111
112     arg_1, arg_2, arg_3 = None, None, None
113
114     for opt, arg in opts:
115         if opt in ['-p']:
116             arg_1 = arg
117         elif opt in ['-e']:
118             arg_2 = arg
119         elif opt in ['-s']:
120             arg_3 = arg
121
122     return arg_1, arg_2, arg_3
123 #-----
124 def usage():
125     sys.exit('Usage:\n'
126             '\t-p: result files directory pathway\n'
127             '\t-e: result files ending\n'
128             '\t-s: species names in the files, separated by commas')
129 #-----
130 if __name__ == '__main__':
131     main()

```

### 8.3 "codeml\_table\_filter.R": corrects p-values using FDR and compile genes for which the bestfit model was M0, bfree or bfree (Topic 5.3)

For reference of data manipulation in R for Chom (codeml\_table\_filter.R):

```

1  setwd("/home/blowflies/CDS_evolution/3-CodeML")
2
3  ##### Chom branch models #####
4
5  # Read table
6  chom_branch <- read.table("Chom_branch_raw_table.txt", header = T)
7
8  # Correct pvalue and get significant results (null = M0)
9  # Does the two omega model (b_free) explain the data better than the one omega
10 # model (M0)?
11 pv_M0 <- chom_branch$pv_free_M0
12 pv_M0_adj <- p.adjust(pv_M0, method = "fdr")
13
14 chom_branch_M0_corrected <- cbind(chom_branch[,1:7], pv_M0_adj, chom_branch[,
15 colnames(chom_branch_M0_corrected)[9] <- colnames(chom_branch)[8]
16
17 chom_branch_M0_significant <- subset(chom_branch_M0_corrected,
18                                     chom_branch_M0_corrected$pv_M0_adj<0.05)
19
20 Chom_bestfit_M0 <- subset(chom_branch_M0_corrected,
21                           chom_branch_M0_corrected$pv_M0_adj>0.05)
22
23 write.table(Chom_bestfit_M0, "Chom_bestfit_M0.txt",
24             sep = "\t", row.names = F, quote = F)
25
26 # Correct pvalue and get significant results (null = b_neut)
27 # Does the free omega in foreground branch model (b_free) explain the data
28 # better than the foreground omega = 1 model (b_neut)?
29 # We only consider the genes that are significant in the first comparison
30 # (if M0 were the best fit, there would be no reason to run this second test)
31
32 pv_neut <- chom_branch_M0_significant$pv_free_neut
33 pv_neut_adj <- p.adjust(pv_neut, method = "fdr")
34
35 chom_branch_neut_corrected <- cbind(chom_branch_M0_significant, pv_neut_adj)
36
37 chom_branch_neut_significant <- subset(chom_branch_neut_corrected,
38                                       chom_branch_neut_corrected$pv_neut_adj<0.05)
39
40 write.table(chom_branch_neut_significant, "Chom_bestfit_free.txt",
41             sep = "\t", row.names = F, quote = F)
42
43 Chom_bestfit_neut <- subset(chom_branch_neut_corrected,
44                             chom_branch_neut_corrected$pv_neut_adj>0.05)
45
46 write.table(Chom_bestfit_neut, "Chom_bestfit_neut.txt",
47             sep = "\t", row.names = F, quote = F)
48
49 # Check if any gene does not have a significant omega
50
51 initial_ids <- chom_branch$Gene
52 significant_ids_M0 <- chom_branch_M0_significant$Gene
53 bestfit_M0_ids <- Chom_bestfit_M0$Gene
54 significant_neut_ids <- chom_branch_neut_significant$Gene
55 bestfit_neut_ids <- Chom_bestfit_neut$Gene
56
57 length(initial_ids)
58 length(bestfit_M0_ids)
59 length(bestfit_neut_ids)
60 length(significant_neut_ids)
61
62 all_ids_combined <- unique(c(Chom_bestfit_M0$Gene, chom_branch_neut_significant$Gene))
63
64 missing_ids <- setdiff(initial_ids, all_ids_combined)
65 print(missing_ids)

```

#### 8.4 "bind-tables\_go-omegas.R": binds tables of GOs and the omegas of their respective associated gene (Topic 6.2) (by ChatGPT)

```

#!/usr/bin/env Rscript

# Usage: Rscript bind-tables_go-omegas.R {/path/to/working/directory} {species} {category}

# Obtain the command line arguments
args <- commandArgs(trailingOnly = TRUE)
if (length(args) < 3) {
  stop("Usage: script.R <path> <species> <category>")
}
work_dir <- args[1]
species <- args[2]
category <- args[3]

# Define working directory
setwd(work_dir)
library(dplyr)

```

```

# Define paths to the temporary directories based on species and category
dir_path <- paste0(species, "_", category, "_tmp/")

# Obtain a list with all files names that finish with .tmp
file_paths <- list.files(path = dir_path, pattern = "\\\\.tmp$", full.names = TRUE)

# Function to read the tables
read_tables <- function(file_paths) {
  lapply(file_paths, read.table, header = TRUE)
}

# Read tables
tables <- read_tables(file_paths)

# Find the maximum number of rows in each table
max_rows <- max(sapply(tables, nrow))

# Function to fill tables with NAs until reaching the maximum number of rows
fill_with_na <- function(df, max_rows) {
  if (nrow(df) < max_rows) df[(nrow(df) + 1):max_rows, ] <- NA
  return(df)
}

# Fill all tables with NA so they have the same number of rows
tables_filled <- lapply(tables, fill_with_na, max_rows)

# Bind tables side by side
merged_df <- bind_cols(tables_filled)

# Write the output file
output_file <- paste0(species, "_", category, "_go-omegas.tsv")
write.table(merged_df, file = output_file, sep = "\t", row.names = FALSE, quote = FALSE)

```

### 8.5 "MW\_test\_multipleX.R": performs Mann-Whitney U test with GOs that have at least X genes associated to them (Topic 6.2) (by ChatGPT)



```

1  #!/usr/bin/env Rscript
2
3  # Load necessary libraries
4  library(dplyr)
5  library(ggplot2)
6
7  # Define arguments
8  args <- commandArgs(trailingOnly = TRUE)
9  if (length(args) < 4) {
10     stop("Usage: script.R <path> <species> <category> <X1,X2,...,XN>")
11 }
12 work_dir <- args[1]
13 species <- args[2]
14 category <- args[3]
15 X_values <- as.numeric(strsplit(args[4], ",")[1])
16
17 # Set working directory
18 setwd(work_dir)
19
20 # Define input files
21 input_file <- paste0(species, "_", category, "_go-omegas.tsv")
22 fixed_column_file <- paste0(species, "_short-infos.tsv")
23
24 # Read the main table
25 df <- read.table(input_file, header = TRUE, sep = "\t")
26
27 # Read the table containing the fixed column
28 fixed_df <- read.table(fixed_column_file, header = TRUE, sep = "\t")
29
30 # Assuming the fixed column is the third column
31 fixed_column <- fixed_df[, 3]
32
33 # Function to perform Mann-Whitney U test and select top 15 headers
34 perform_tests <- function(X) {
35     # Filter columns with at least X non-NA values
36     filtered_df <- df %>% select_if(~ sum(!is.na(.)) >= X)
37
38     # Perform Mann-Whitney U test for each selected column and store p-values
39     p_values <- sapply(filtered_df, function(column) {
40         wilcox.test(column, fixed_column, na.rm = TRUE)$p.value
41     })
42
43     # Apply FDR correction to p-values
44     p_values_adj <- p.adjust(p_values, method = "fdr")
45
46     # Calculate medians of each column
47     medians <- sapply(filtered_df, median, na.rm = TRUE)
48
49     # Calculate the median of the fixed column
50     fixed_median <- median(fixed_column, na.rm = TRUE)
51
52     # Headers of columns with significantly higher medians
53     higher_significant_headers <- colnames(filtered_df)[p_values_adj < 0.05 &
54     higher_significant_headers <- higher_significant_headers[order(medians[high
55
56     # Headers of columns with significantly lower medians
57     lower_significant_headers <- colnames(filtered_df)[p_values_adj < 0.05 &
58     lower_significant_headers <- lower_significant_headers[order(medians[low
59
60     return(list(
61         higher_significant_headers = higher_significant_headers,
62         lower_significant_headers = lower_significant_headers
63     ))
64 }
65
66 # Execute tests for each X value
67 results <- lapply(X_values, perform_tests)
68
69 # Extract results
70 higher_significant_headers_list <- lapply(results, function(result) result$high
71 lower_significant_headers_list <- lapply(results, function(result) result$low
72
73 # Create tables for top 15 higher and lower headers
74 top15_higher_table <- data.frame(matrix(nrow = 15, ncol = length(X_values)))
75 top15_lower_table <- data.frame(matrix(nrow = 15, ncol = length(X_values)))
76 colnames(top15_higher_table) <- paste0("X", X_values)
77 colnames(top15_lower_table) <- paste0("X", X_values)
78
79 for (i in 1:length(X_values)) {
80     top15_higher_table[, i] <- higher_significant_headers_list[[i]]
81     top15_lower_table[, i] <- lower_significant_headers_list[[i]]
82 }
83
84 # Save the tables
85 write.table(top15_higher_table, file = paste0(species, "_", category, "_top15
86 write.table(top15_lower_table, file = paste0(species, "_", category, "_top15
87
88 # Function to perform Mann-Whitney U test and count the number of significant

```

```

89 perform_test_counts <- function(X) {
90   # Filter columns with at least X non-NA values
91   filtered_df <- df %>% select_if(~ sum(!is.na(.)) >= X)
92
93   # Perform Mann-Whitney U test for each selected column and store p-values
94   p_values <- sapply(filtered_df, function(column) {
95     wilcox.test(column, fixed_column, na.rm = TRUE)$p.value
96   })
97
98   # Apply FDR correction to p-values
99   p_values_adj <- p.adjust(p_values, method = "fdr")
100
101   # Count the number of significant tests (adjusted p-value < 0.05)
102   num_significant <- sum(p_values_adj < 0.05)
103
104   return(num_significant)
105 }
106
107 # Execute tests for each X value and count the significant tests
108 num_significant <- sapply(X_values, perform_test_counts)
109
110 # Create a data frame for the plot
111 df_plot <- data.frame(X = X_values, num_significant = num_significant)
112
113 # Generate the plot
114 plot_file <- paste0(species, "_", category, "_num_significant_plot.png")
115 ggplot(df_plot, aes(x = X, y = num_significant)) +
116   geom_line() +
117   geom_point() +
118   labs(title = "Number of Significant Tests vs. X", x = "X", y = "Number of Significant Tests") +
119   theme_minimal() +
120   ggsave(plot_file)
121
122 print(paste("Graph saved to", plot_file))

```

**8.6 "MW\_test\_table.R": creates a table that compiles info necessary to plot a multiple boxplot graph (by ChatGPT)**





```

1  #!/usr/bin/env Rscript
2
3  # Load necessary libraries
4  library(dplyr)
5
6  # Define arguments
7  args <- commandArgs(trailingOnly = TRUE)
8  if (length(args) < 4) {
9      stop("Usage: script.R <path> <species> <X_value> <category1,category2,..>")
10 }
11 work_dir <- args[1]
12 species <- args[2]
13 X_value <- as.numeric(args[3])
14 categories <- strsplit(args[4], ",")[1]
15
16 # Set working directory
17 setwd(work_dir)
18
19 # Define the input file for the fixed column
20 fixed_column_file <- paste0(species, "_short-infos.tsv")
21
22 # Read the table containing the fixed column
23 fixed_df <- read.table(fixed_column_file, header = TRUE, sep = "\t")
24
25 # Assuming the fixed column is the third column
26 fixed_column <- fixed_df[, 3]
27
28 # Function to perform Mann-Whitney U test and select top 15 headers
29 perform_tests <- function(df, X) {
30     # Filter columns with at least X non-NA values
31     filtered_df <- df %>% select_if(~ sum(!is.na(.)) >= X)
32
33     # Perform Mann-Whitney U test for each selected column and store p-value
34     p_values <- sapply(filtered_df, function(column) {
35         wilcox.test(column, fixed_column, na.rm = TRUE)$p.value
36     })
37
38     # Apply FDR correction to p-values
39     p_values_adj <- p.adjust(p_values, method = "fdr")
40
41     # Calculate medians of each column
42     medians <- sapply(filtered_df, median, na.rm = TRUE)
43
44     # Calculate the median of the fixed column
45     fixed_median <- median(fixed_column, na.rm = TRUE)
46
47     # Headers of columns with significantly higher medians
48     higher_significant_headers <- colnames(filtered_df)[p_values_adj < 0.05 & medians > fixed_median]
49     higher_significant_headers <- higher_significant_headers[order(medians[higher_significant_headers], decreasing = TRUE)]
50
51     # Headers of columns with significantly lower medians
52     lower_significant_headers <- colnames(filtered_df)[p_values_adj < 0.05 & medians < fixed_median]
53     lower_significant_headers <- lower_significant_headers[order(medians[lower_significant_headers], increasing = TRUE)]
54
55     return(list(
56         higher_significant_headers = higher_significant_headers,
57         lower_significant_headers = lower_significant_headers,
58         filtered_df = filtered_df
59     ))
60 }
61
62 # Initialize lists to store results
63 category_list <- c()
64 header_list <- c()
65 values_list <- c()
66
67 # Process each category
68 for (category in categories) {
69     # Define the input file for the category
70     input_file <- paste0(species, "_", category, "_go-omegas.tsv")
71
72     # Read the table for the category
73     df <- read.table(input_file, header = TRUE, sep = "\t")
74
75     # Perform tests to get significant headers
76     results <- perform_tests(df, X_value)
77
78     # Extract results
79     higher_significant_headers <- results$higher_significant_headers
80     lower_significant_headers <- results$lower_significant_headers
81     filtered_df <- results$filtered_df
82
83     # Collect data for higher significant headers
84     for (header in higher_significant_headers) {
85         if (!is.na(header)) {
86             values <- paste(na.omit(filtered_df[[header]]), collapse = ", ")
87             category_list <- c(category_list, category)
88             header_list <- c(header_list, header)

```

```

89     values_list <- c(values_list, values)
90   }
91 }
92
93 # Collect data for lower significant headers
94 for (header in lower_significant_headers) {
95   if (!is.na(header)) {
96     values <- paste(na.omit(filtered_df[[header]]), collapse = ", ")
97     category_list <- c(category_list, category)
98     header_list <- c(header_list, header)
99     values_list <- c(values_list, values)
100   }
101 }
102 }
103
104 # Collect data for the fixed column
105 fixed_values <- paste(na.omit(fixed_column), collapse = ", ")
106 category_list <- c(category_list, rep("Genome", 1))
107 header_list <- c(header_list, rep("Genome", 1))
108 values_list <- c(values_list, fixed_values)
109
110 # Create the final data frame
111 final_df <- data.frame(Category = category_list, GO_ID = header_list, List_C
112
113 # Save the final table
114 output_file <- paste0(species, "_top30_GOs.tsv")
115 write.table(final_df, file = output_file, sep = "\t", row.names = FALSE, qu

```

### 8.7 "GO\_terms.py": inputs a GO id list and outputs a table with ids and term names (by ChatGPT)

```

1  import requests
2  import pandas as pd
3
4  # Lista de GO IDs
5  go_ids = [
6      "GO:0050912", "GO:0050909", "GO:0050907", "GO:0009593", "GO:0010951", "GO:0010952",
7  ]
8
9  # Função para obter os nomes completos dos termos GO usando a API do QuickGO
10 def get_go_names(go_ids):
11     base_url = "https://www.ebi.ac.uk/QuickGO/services/ontology/go/terms/"
12     headers = {"Accept": "application/json"}
13     go_names = []
14
15     for go_id in go_ids:
16         response = requests.get(base_url + go_id, headers=headers)
17         if response.status_code == 200:
18             data = response.json()
19             go_names.append({"GO ID": go_id, "Term Name": data['results'][0]['name']})
20         else:
21             go_names.append({"GO ID": go_id, "Term Name": "Not found"})
22
23     return go_names
24
25 # Obter os nomes completos dos termos GO
26 go_names = get_go_names(go_ids)
27
28 # Converter a lista de dicionários em um DataFrame
29 df_go_names = pd.DataFrame(go_names)
30
31 # Salvar em um arquivo CSV (opcional)
32 df_go_names.to_csv("go_term_names.csv", index=False)
33
34 # Exibir a tabela
35 print(df_go_names)

```

### 8.8 "GO\_boxplot.R": script to create figures like the one from Scott et al. (2020), for the rate of evolution of the genes associated to specific GOs



```

1 library(tidyr)
2 library(dplyr)
3 library(ggplot2)
4 setwd("C:/Users/Pedro/OneDrive/Área de Trabalho/GO_figure")
5
6 ##### Reading the files #####
7 Chom_genome <- read.table("Chom_genome_omegas.tsv", sep="\t", header=T)
8 Chom_top30 <- read.table("Chom_top30_GOs.tsv", sep="\t", header=F)
9
10 Pazu_genome <- read.table("Pazu_genome_omegas.tsv", sep="\t", header=T)
11 Pazu_top30 <- read.table("Pazu_top30_GOs.tsv", sep="\t", header=F)
12
13 ##### Fixing the top30 tables #####
14 # Chom
15 colnames(Chom_top30)[1] <- "Category" # fixing the header
16 colnames(Chom_top30)[2] <- "GO_id"
17 colnames(Chom_top30)[3] <- "GO_term"
18
19 Chom_top30_long <- gather(Chom_top30, Number, Omega, V4:V580) # wide to long
20
21 Chom_top30_long <- na.omit(Chom_top30_long) # remove NA, because there differ
22 # quantity of genes with each GO term
23
24 # Pazu
25 colnames(Pazu_top30)[1] <- "Category" # fixing the header
26 colnames(Pazu_top30)[2] <- "GO_id"
27 colnames(Pazu_top30)[3] <- "GO_term"
28
29 Pazu_top30_long <- gather(Pazu_top30, Number, Omega, V4:V537) # wide to long
30
31 Pazu_top30_long <- na.omit(Pazu_top30_long) # remove NA, because there differ
32 # quantity of genes with each GO term
33
34 ##### Fixing genome file and merging #####
35 # Chom
36 G <- rep("Genome", nrow(Chom_genome)) # make the genome table equal to the G
37 Chom_genome <- cbind(G,G,G,G,Chom_genome)
38 colnames(Chom_genome)[1] <- "Category"
39 colnames(Chom_genome)[2] <- "GO_id"
40 colnames(Chom_genome)[3] <- "GO_term"
41 colnames(Chom_genome)[4] <- "Number"
42
43 Chom_final <- rbind(Chom_top30_long,Chom_genome) # merge everything
44
45 Chom_final$Category <- factor(Chom_final$Category,
46 levels = c("BP", "MF", "CC", "Genome")) # fix
47 Chom_final$GO_term <- as.factor(Chom_final$GO_term)
48 Chom_final$Omega <- as.numeric(Chom_final$Omega)
49
50 Chom_final_medians <- Chom_final %>% # organize GO terms by omega mean
51 group_by(GO_term) %>%
52 summarize(median_Omega = median(Omega, na.rm = TRUE))
53
54 Chom_final <- Chom_final %>%
55 mutate(GO_term = factor(GO_term, levels = Chom_final_medians$GO_term[order
56
57 # Pazu
58 G <- rep("Genome", nrow(Pazu_genome)) # make the genome table equal to the G
59 Pazu_genome <- cbind(G,G,G,G,Pazu_genome)
60 colnames(Pazu_genome)[1] <- "Category"
61 colnames(Pazu_genome)[2] <- "GO_id"
62 colnames(Pazu_genome)[3] <- "GO_term"
63 colnames(Pazu_genome)[4] <- "Number"
64
65 Pazu_final <- rbind(Pazu_top30_long,Pazu_genome) # merge everything
66
67 Pazu_final$Category <- factor(Pazu_final$Category,
68 levels = c("BP", "MF", "CC", "Genome")) # fix
69 Pazu_final$GO_term <- as.factor(Pazu_final$GO_term)
70 Pazu_final$Omega <- as.numeric(Pazu_final$Omega)
71
72 Pazu_final_medians <- Pazu_final %>% # organize GO terms by omega mean
73 group_by(GO_term) %>%
74 summarize(median_Omega = median(Omega, na.rm = TRUE))
75
76 Pazu_final <- Pazu_final %>%
77 mutate(GO_term = factor(GO_term, levels = Pazu_final_medians$GO_term[order
78
79 ##### Figure #####
80 # Chom
81 chom_genome_median = median(Chom_genome$Omega)
82
83 svg("Chom_GO_fig.svg", width=15, height=15)
84 ggplot(Chom_final, aes(x=Omega,y=GO_term,fill=Category)) +
85 geom_boxplot() +
86 facet_grid(Category ~ ., scales = "free") +
87 labs(title = "Cochliomyia hominivorax", x="dN/dS ratio", y="GO terms") +
88 theme(plot.title = element_text(face = "italic")) +

```

```

89     geom_vline(xintercept=chom_genome_median, linetype="dashed", color="red")
90     dev.off()
91
92     png("Chom_GO_fig.png", width=1000, height=1000)
93     ggplot(Chom_final, aes(x=Omega,y=GO_term,fill=Category)) +
94       geom_boxplot () +
95       facet_grid(Category ~ ., scales= "free") +
96       labs(title = "Cochliomyia hominivorax", x="dN/dS ratio", y="GO terms") +
97       theme(plot.title = element_text(face = "italic")) +
98       geom_vline(xintercept=chom_genome_median, linetype="dashed", color="red")
99     dev.off()
100
101     # Pazu
102     pazu_genome_median = median(Pazu_genome$Omega)
103
104     svg("Pazu_GO_fig.svg", width=15, height=15)
105     ggplot(Pazu_final, aes(x=Omega,y=GO_term,fill=Category)) +
106       geom_boxplot () +
107       facet_grid(Category ~ ., scales= "free") +
108       theme(plot.title = element_text(face = "italic")) +
109       labs(title = "Protocalliphora azurea", x="dN/dS ratio", y="GO terms") +
110       geom_vline(xintercept=pazu_genome_median, linetype="dashed", color="red")
111     dev.off()
112
113     png("Pazu_GO_fig.png", width=1000, height=1000)
114     ggplot(Pazu_final, aes(x=Omega,y=GO_term,fill=Category)) +
115       geom_boxplot () +
116       facet_grid(Category ~ ., scales= "free") +
117       labs(title = "Protocalliphora azurea", x="dN/dS ratio", y="GO terms") +
118       theme(plot.title = element_text(face = "italic")) +
119       geom_vline(xintercept=pazu_genome_median, linetype="dashed", color="red")
120     dev.off()

```

## 8.9 "get\_ids\_from\_cluster\_list.py"

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3
4  '''
5  PROGRAM: get_sequences_from_id_list.py
6
7  DESCRIPTION: Reads a file with a list of ortholog clusters identifiers and
8               recovers the sequences belonging to a given species from these
9               clusters fasta files.
10
11  AUTHOR: Pedro Mariano Martins                São Paulo, Brazil, 09.08.2020
12  '''
13
14  import sys
15  import getopt
16  import glob
17  #-----
18  def main():
19      ids, species, cluster_pwd, cluster_end = arg()
20
21      if ids == None or species == None or cluster_pwd == None or cluster_end == None:
22          usage ()
23
24      clusters = glob.glob(f'{cluster_pwd}/{cluster_end}')
25
26      #making a list with the ids
27      idents = open(ids, 'r', encoding='utf-8')
28      id_list = []
29
30      for line in idents:
31          line = line.strip()
32          id_list += [line]
33
34      idents.close()
35
36      #finding clusters and sequences, and saving them in an output file
37      out = open(f'{species}_table_from_clusters.txt', 'w', encoding='utf-8')
38
39      for cluster in clusters:
40          NA = 0
41          cluster_name = cluster.split(f'{cluster_pwd}/')[1]
42          for i in id_list:
43              if i in cluster:
44                  fasta = open(cluster, 'r', encoding='utf-8')
45                  for line in fasta:
46                      if species in line:
47                          line = line.strip()
48                          out.write(f'{cluster_name}\t{line}\n')
49                      NA = 1
50                  fasta.close()
51              if NA == 0:
52                  out.write(f'{cluster_name}\tNA\n')
53
54      out.close()
55  #-----
56  def arg():
57      argv = sys.argv[1:]
58      try:
59          opts, args = getopt.getopt(argv, 'i:s:c:e:')
60
61      except:
62          usage()
63
64      arg_1, arg_2, arg_3, arg_4 = None, None, None, None
65
66      for opt, arg in opts:
67          if opt in ['-i']:
68              arg_1 = arg
69          elif opt in ['-s']:
70              arg_2 = arg
71          elif opt in ['-c']:
72              arg_3 = arg
73          elif opt in ['-e']:
74              arg_4 = arg
75
76      return arg_1, arg_2, arg_3, arg_4
77  #-----
78  def usage():
79      sys.exit('Usage:\n'
80              '\t-i: list of cluster names\n'
81              '\t-s: species name\n'
82              '\t-c: cluster files pathway\n'
83              '\t-e: cluster files ending')
84  #-----
85  if __name__ == '__main__':

```

