

Osedax paper analyses

Chapter 3

k-mer analyses

kraken2

```
1 module load anaconda3
2 conda create -n kraken2_env2 -c bioconda kraken2
3 conda install -c conda-forge libiconv
```

- kraken2 v 2.1.0

krakentools

```
1 module load anaconda3
2 conda activate kraken2_env
3 conda install -c bioconda krakentools
```

- krakentools v 0.1

Kraken & KrakenTools

kraken2_pacbio_bbmap_standard_db_lessCores_v1.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/
3 #$ -o /data/scratch/btx654/
4 #$ -j y
5 #$ -pe smp 6
6 #$ -l h_vmem=10G
7 #$ -l h_rt=240:0:0
8 #$ -l highmem
9
10 species=$1
```

```
11 kraken_database=/data/SBCS-MartinDuranLab/03-Giacomo/db/kraken2_
    db_standard_Oct2020
12 pacbio_corrected="$species"_pacbio_corrected_bbmap.fasta
13 output_pacbio="$species"_pacbio_bbmap_output.kraken
14
15 module load anaconda3
16 conda activate kraken2_env
17
18 cd /data/scratch/btx654/btx604-scratch/$species/kraken2/pacbio_i
    llumina_2
19
20 kraken2 --threads 6 --output $output_pacbio --report report_krak
    en2_pacbio_bbmap --db $kraken_database $pacbio_corrected
```

Kraken tools extract oasisia

```
1 module load anaconda3
2 conda activate kraken2_env
3 extract_kraken_reads.py -k oasisia_pacbio_bbmap_output.kraken -r
    report_kraken2_pacbio_bbmap -s oasisia_pacbio_corrected_bbmap.fa
    sta --taxid 2 --exclude --include-children -o oasisia_nonBacteri
    a_pacbio_corrected_bbmap.fasta
```

Illumina nonBacterial

mapping_illumina_nonBacterial_definitive.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654
3 #$ -j y
4 #$ -o /data/scratch/btx654
5 #$ -pe smp 20
6 #$ -l h_vmem=34G
7 #$ -l h_rt=72:0:0
8 #$ -l highmem
```

```
9
10 species=$1
11 #bwa index step1 variables
12 pacbio_corrected_nonBacteria="$species"_nonBacteria_pacbio_corre
    cted_bbmap.fasta
13 bwa_prefix="$species"_nonBacteria_pacbio_corrected_bbmap.fasta
14 #bwa mem step1 variables
15 R1_cleaned="$species"_R1_cleaned.fastq.gz
16 R2_cleaned="$species"_R2_cleaned.fastq.gz
17 alignment_sam="$species"_alignment_HIGHmem.sam
18 #samtools view step1 variables
19 alignment_bam="$1"_alignment_HIGHmem.bam
20 #samtools sort index step1 variables
21 alignment_sorted="$species"_sorted_HIGHmem.bam
22 alignment_mapped="$species"_sorted_mapped_HIGHmem.bam
23 alignment_mapped_sorted="$species"_sorted_mapped_sorted_HIGHmem.
    bam
24 R1_mapped="$species"_R1_mapped_HIGHmem.fastq
25 R2_mapped="$species"_R2_mapped_HIGHmem.fastq
26
27 module load bwa
28 module load samtools/1.9
29
30 cd /data/scratch/btx654/btx604-scratch/$species/kraken2/pacbio_i
    llumina
31
32 echo 'BWA INDEX STEP1_-----
    -----'
33 if [ -e /data/scratch/btx654/btx604-scratch/$species/kraken2/pac
    bio_illumina/*.ann ]
34 then
```

```

35     echo "/data/scratch/btx654/btx604-scratch/$species/kraken2/pac
bio_illumina/*.ann found."
36 else
37     bwa index -p $bwa_prefix -a bwtsv $pacbio_corrected_nonBacteri
a
38 fi
39
40 echo 'BWA MEM STEP1_____
_____
41 bwa mem -t 20 -M $pacbio_corrected_nonBacteria $R1_cleaned $R2_c
leaned > $alignment_sam
42
43 samtools view -@ 20 -S -b -h $alignment_sam -o $alignment_bam
44 samtools sort -@ 20 $alignment_bam -o $alignment_sorted
45 samtools index $alignment_sorted
46 samtools view -@ 20 -b -F 4 $alignment_sorted > $alignment_mappe
d
47 samtools sort -@ 20 -n $alignment_mapped -o $alignment_mapped_so
rted
48 samtools fastq -1 illumina_R1_mapped_samtools.fq -2 illumina_R2_
mapped_samtools.fq -n $alignment_mapped_sorted

```

KAT

First let's set up a directory containing illumina reads for Riftia and the nonBacterial illumina reads obtained by mapping illumina reads on the non-bacterial pacbio generated by kraken2 (29/10/20). For Riftia the illumina needs to be cleaned using fastp

/data/home/btx654/scripts/kmer_analyses/Dec2020/fastp_riftia_v1.sh

```

1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/

```

```
4  #$ -j y
5  #$ -pe smp 4
6  #$ -l h_vmem=5G
7  #$ -l h_rt=6:0:0
8
9  species=riftia
10 R1=RP_6_12_19_S110_L002_R1_001.fastq.gz
11 R2=RP_6_12_19_S110_L002_R2_001.fastq.gz
12 R1_cleaned="$species"_R1.fq.gz
13 R2_cleaned="$species"_R2.fq.gz
14
15 cd /data/scratch/btx654/btx604-scratch/$species/kmer_Dec2020/
16
17 module load anaconda3
18 conda activate fastp
19
20 fastp -i $R1 -I $R2 -o $R1_cleaned -O $R2_cleaned -w 4
21
22 gzip -d $R1_cleaned
23 gzip -d $R2_cleaned
```

fastp_osedax_oasisia_universal_v1.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 4
6  #$ -l h_vmem=5G
7  #$ -l h_rt=6:0:0
8
9  species=$1
```

```
10 R1=illumina_R1_nonBacteria_kraken2.fq
11 R2=illumina_R2_nonBacteria_kraken2.fq
12 R1_cleaned="$species"_R1.fq
13 R2_cleaned="$species"_R2.fq
14
15 cd /data/scratch/btx654/btx604-scratch/$species/kmer_Dec2020/
16
17 module load anaconda3
18 conda activate fastp
19
20 fastp -i $R1 -I $R2 -o $R1_cleaned -O $R2_cleaned -w 4
```

After this we want to have our folders:

```
/data/scratch/btx654/btx604-scratch/$species/kmer_Dec2020/
```

containing illumina data named as "\$species"_R1.fq and "\$species"_R2.fq (uncompressed)

- for Riftia this 2 files are the fastp cleaned illumina data we had from 2019 (not the august 2020 data)
- for Osedax and Oasisia this 2 files are the nonBacterial illumina reads obtained by mapping illumina reads on the non-bacterial pacbio generated by kraken2 (29/10/20). Which have been cleaned with the fastp cleaning step

WILL BASE ALL THESE ANALYSES ON 21-MER!!!!

kat_gcp_universal_v1.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/
3 #$ -o /data/scratch/btx654/
4 #$ -j y
5 #$ -pe smp 20
6 #$ -l h_vmem=10G
7 #$ -l h_rt=24:0:0
```

```
8  #$ -l highmem
9
10 species=$1
11 R1="$species"_R1.fq
12 R2="$species"_R2.fq
13
14 echo "Working on "$species
15
16 cd /data/scratch/btx654/btx604-scratch/$species/kmer_Dec2020/
17 mkdir -p kat
18 cd kat
19
20 module load anaconda3
21 conda activate KAT
22
23 kat gcp -m 21 -p pdf -v -t 20 ../$R1 ../$R2
```

kat_hist_universal_v1.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 20
6  #$ -l h_vmem=10G
7  #$ -l h_rt=72:0:0
8  #$ -l highmem
9
10
11 species=$1
12 R1="$species"_R1.fq
13 R2="$species"_R2.fq
```

```
14
15 echo "Working on "$species
16
17 cd /data/scratch/btx654/btx604-scratch/$species/kmer_Dec2020/kat
18
19 module load anaconda3
20 conda activate KAT
21
22 kat hist -m 21 -t 20 -o "kat_21mer_illumina.hist" ../$R1 ../$R2
```

Before the next step cp the purged genome into the working directory and name it "\$species".fa

kat_comp_universal_v1.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/
3 #$ -o /data/scratch/btx654/
4 #$ -j y
5 #$ -pe smp 20
6 #$ -l h_vmem=10G
7 #$ -l h_rt=24:0:0
8 #$ -l highmem
9
10 species=$1
11 purged="$species".fa
12
13 echo "Working on "$species
14
15 module load anaconda3
16 conda activate KAT
17
```



```
18 cd /data/scratch/btx654/btx604-scratch/$species/kmer_Dec2020/kat
19
20 kat comp -m 21 -p pdf -o 21mer_vs_assembly -v -t 20 '../*_R1.fq
  ../*_R2.fq' ../$purged
```

mercury_universal.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 10
6  #$ -l h_vmem=20G
7  #$ -l h_rt=100:0:0
8  #$ -l highmem
9
10 species=$1
11 purged="$species".fa
12 R1="$species"_R1.fq
13 R2="$species"_R2.fq
14 meryl_R1="$species"_R1.meryl
15 meryl_R2="$species"_R2.meryl
16 meryl_final="$species".meryl
17 mercury_output="$species"_mercury_nonBacteria
18
19 echo "Working on "$species
20
21 cd /data/scratch/btx654/btx604-scratch/$species/kmer_Dec2020/
22 mkdir -p mercury
23 cd mercury
24
```

```
25 module load anaconda3
26 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
3/mercury_env
27
28 meryl k=21 threads=10 count output $meryl_R1 ../$R1
29 meryl k=21 threads=10 count output $meryl_R2 ../$R2
30
31 meryl union-sum output $meryl_final *_R*.meryl
32
33 mercury.sh $meryl_final ../$purged $mercury_output > Mercury.log
34
35 #spectra_cn="$species"_mercury_nonBacteria."$species".spectra-c
36 n.hist
37 #only_hist="$species"_mercury_nonBacteria.dist_only.hist
38
39 #output_plot="$species".spectra-cn
40
41 #Rscript $MERQUERY/plot/plot_spectra_cn.R -f $spectra_cn -o $outp
42 ut_plot -z $only_his -m (kmer_multiplicity) and -n (Count)
```

Genome size estimation -GenomeScope

Riftia - <http://qb.cshl.edu/genomescope/genomescope2.0/analysis.php?code=XcAsfjFyAhGcklnVGStm>

Osedax - <http://qb.cshl.edu/genomescope/genomescope2.0/analysis.php?code=aUtpCdYPBD5MXkT1h0cb>

Oasisia - <http://qb.cshl.edu/genomescope/genomescope2.0/analysis.php?code=5hnEsmQKmAf6j8u65ub7>

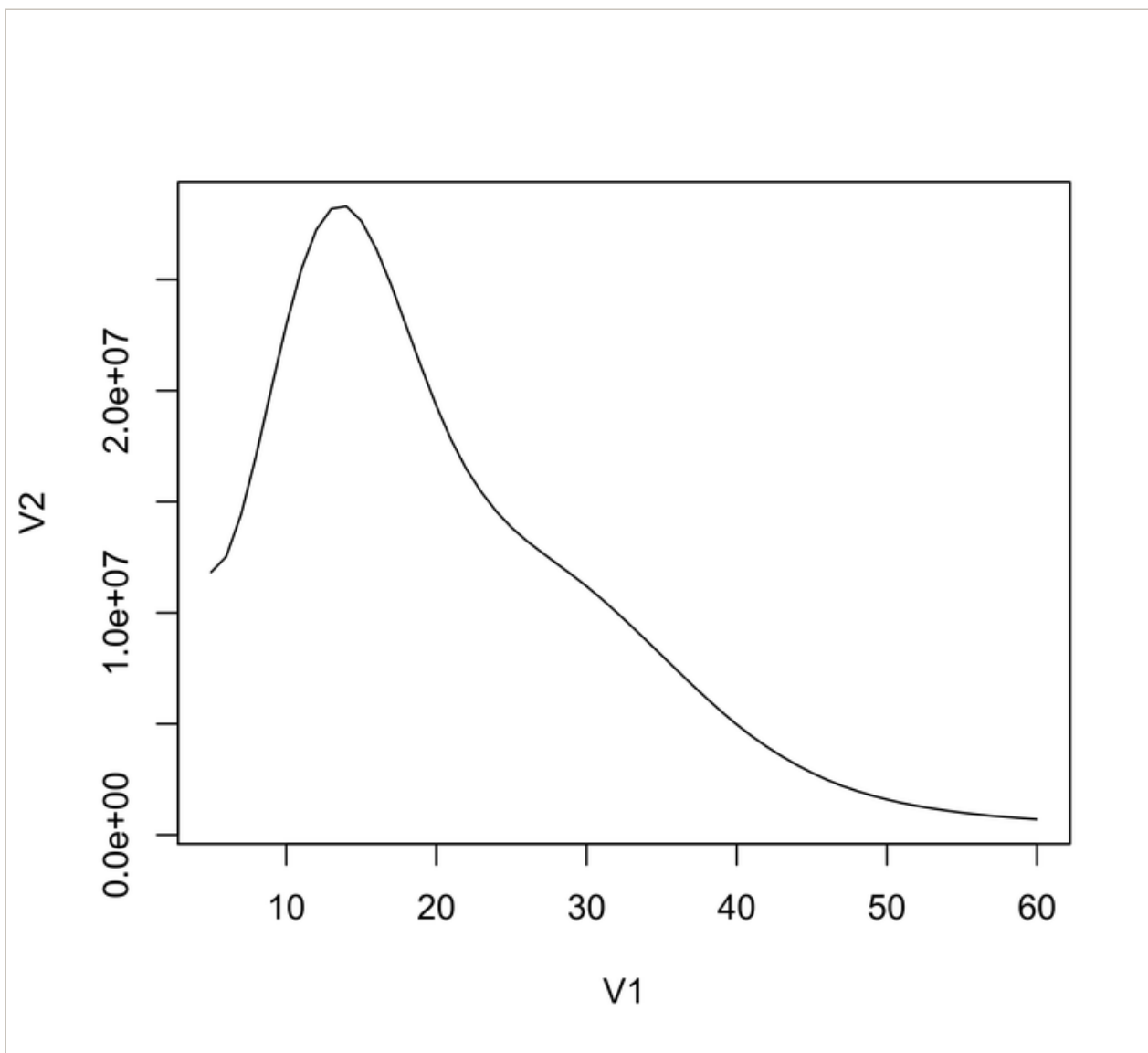
Genome size estimation - [tutorial](#)

R based calculations

Using the hist file produced with the script

(kat_hist_universal_kraken_standard_illumina_v1.sh 28/10/20) and cropping out of the plot the first area of erroneous kmers as suggested by the tutorial

```
1 plot(oasisia_21mer[5:60,], type="l")
2 # points(oasisia_21mer[5:60,])
3 # export as a pdf 5 x 5.5 inches in landscape orientation
```

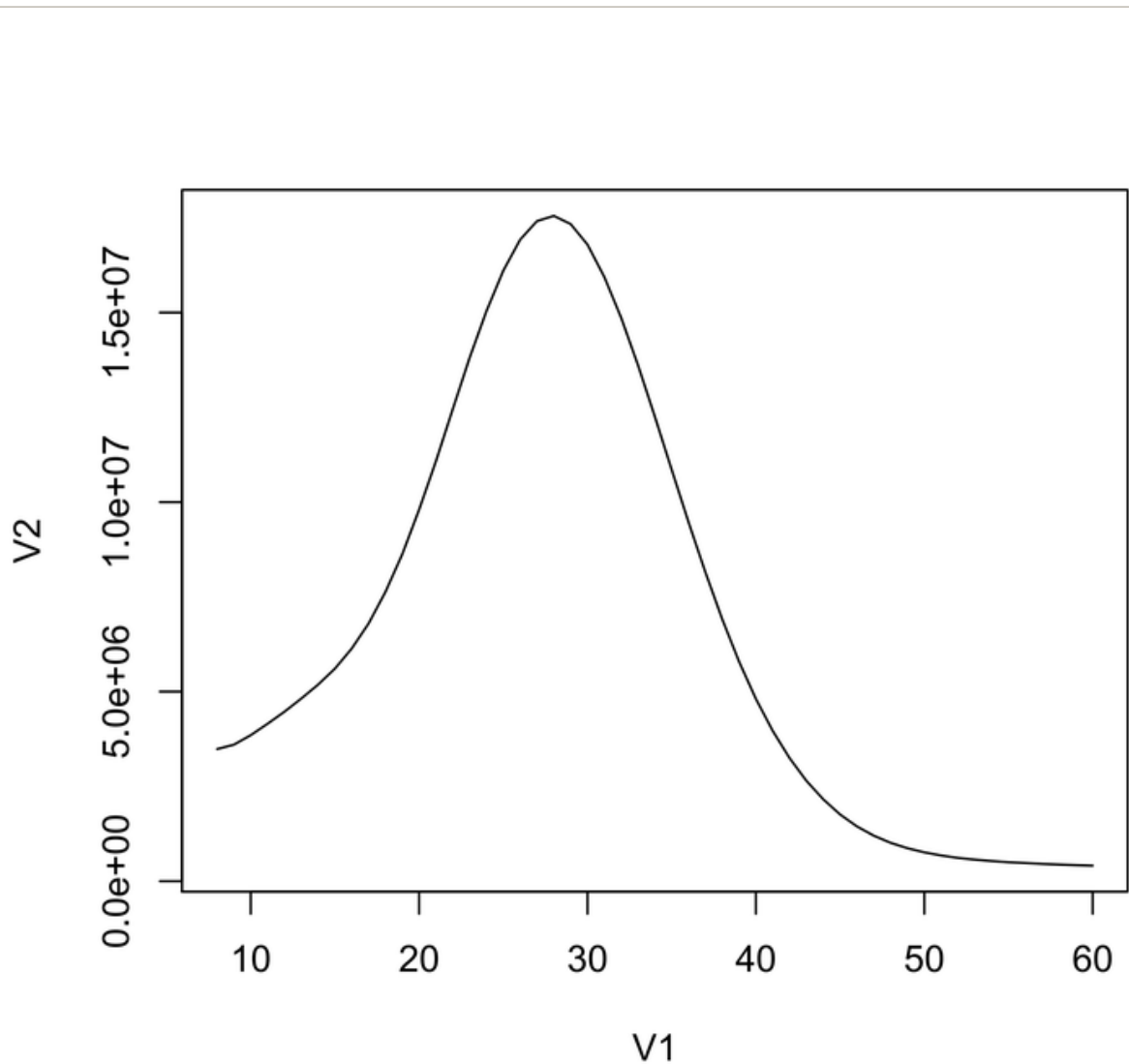


PDF oasisia_21mer_spectra • PDF document

```
1 #here I am assuming I have 10000 data points (which is the length
  # of the hist file)
2 # 28 is the peak for the homozygous and the result is similar to
  # my assembly size
```

```
3 > sum(as.numeric(oasisia_21mer[5:10000,1]*oasisia_21mer[5:10000,
4 [1] 785764774
```

```
1 plot(riftia_21mer[8:60,], type="l")
2 # points(oasisia_21mer[5:60,])
3 # export as a pdf 5 x 5.5 inches in landscape orientation
```



PDF riftia_21mer_spectra • PDF document

```
1 #here I am assuming I have 10000 data points (which is the lengt
```

```
h of the hist file)
```

```
2 # 27 is the peak for the homozygous and the result is similar to  
my assembly size
```

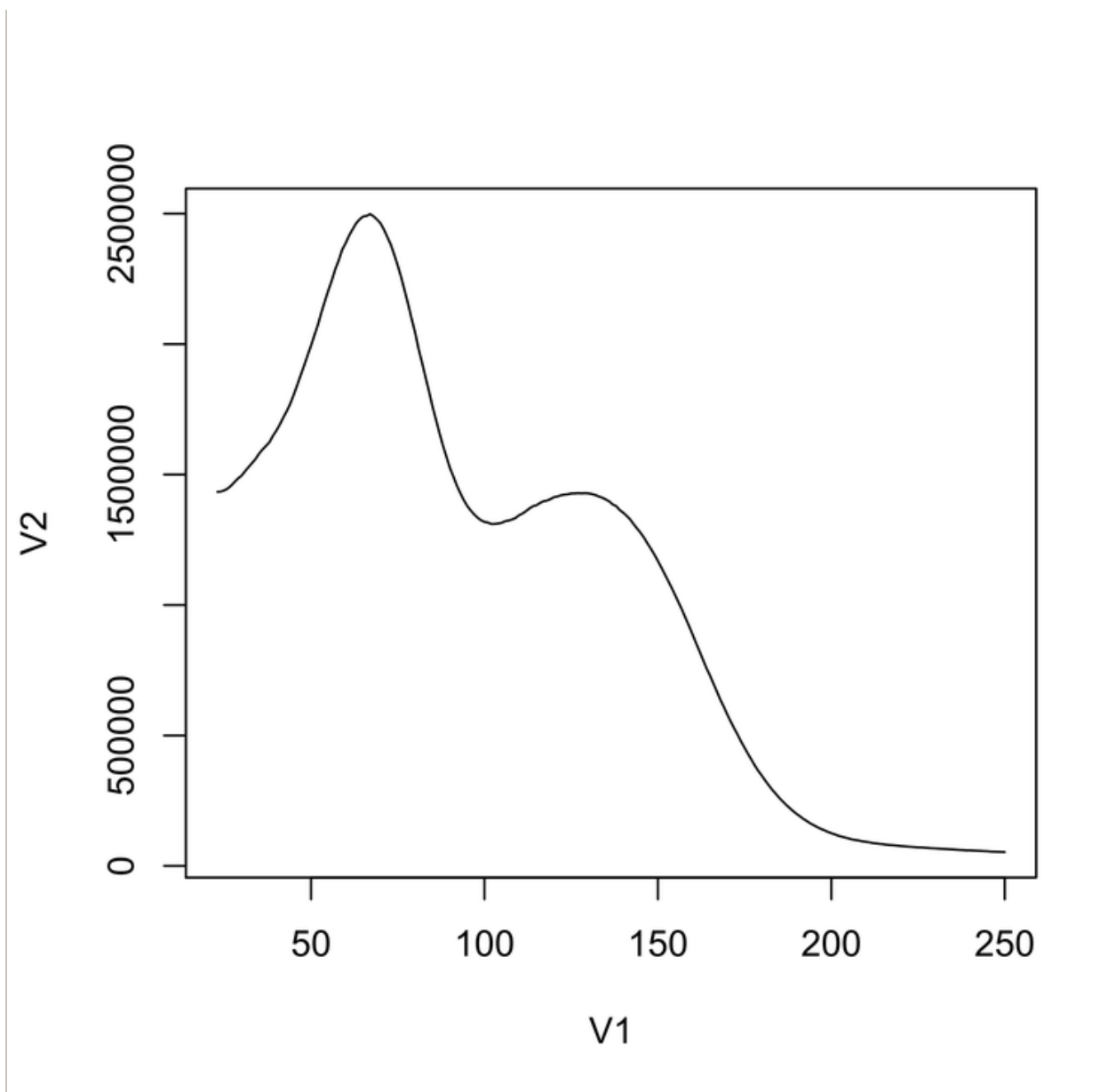
```
3 > sum(as.numeric(riftia_21mer[8:10000,1]*riftia_21mer[8:10000,  
2]))/28
```

```
4 [1] 532856757
```

```
1 plot(osedax_21mer[23:250,], type="l")
```

```
2 # points(oasisia_21mer[5:60,])
```

```
3 # export as a pdf 5 x 5.5 inches in landscape orientation
```



PDF osedax_21mer_spectra • PDF document

```
1 #here I am assuming I have 10000 data points (which is the length of the hist file)
2 # 127 is the peak for the homozygous and the result is similar to my assembly size
3 > sum(as.numeric(osedax_21mer[23:10000,1]*osedax_21mer[23:10000,2]))/127
4 [1] 239329736
```

Missing Busco in Osedax

files I will use:

```
1 /Users/giacomo/Dropbox/11-Siboglinids/00-Data/Osedax/Annotation/
  New_annotation_Dec2020/step7/missing_busco_list.tsv
2 /data/SBCS-MartinDuranLab/03-Giacomo/db/datasets/metazoa_odb10/a
  ncestral
```

busco_universal_v1.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/
3 #$ -o /data/scratch/btx654/
4 #$ -pe smp 4
5 #$ -l h_vmem=20G
6 #$ -l h_rt=48:0:0
7 #$ -j y
8 #$ -l highmem
9
10 species=osedax
11 annotation_gtf="$species".AGAT.noSTOP.filt.noTE.gtf
12 annotation_fa="$species"_annotation.prot.fa
13 species_softmasked="$species"_softmasked.fa
14 output_busco="$species"_busco_annotation
15
16 echo "Working on "$species
17
18 cd /data/scratch/btx654/missing_busco_osedax
19 cp /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax/annotation/N
  ew_annotation_Dec2020/step6/$annotation_gtf ./
20 cp /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax/annotation/s
  oftmasking/$species_softmasked ./
21
```

```
22 module load anaconda3
23 source activate augustus
24
25 gffread -E $annotation_gtf -g $species_softmasked -y $annotation
_fa
26
27 conda deactivate
28 source activate busco_env
29 #export BUSCO_CONFIG_FILE="/data/home/btx654/.conda/envs/busco_e
nv/busco/config/myconfig.ini"
30 #export AUGUSTUS_CONFIG_PATH=/data/SBCS-MartinDuranLab/02-Chema/
src/Augustus/config/
31
32 busco -i osedax_mRNA.fa -m proteins -o output_busco -c 4 -l /dat
a/SBCS-MartinDuranLab/03-Giacomo/db/datasets/metazoa_odb10
```

```
1 sed 's/a//g' missing_busco_list.tsv | sed 's/^/>/' > missing_bus
co_list_OK.txt
2
3 cp /data/SBCS-MartinDuranLab/03-Giacomo/db/datasets/metazoa_odb1
0/ancestral ./
4 mv ancestral ancestral.fa
5
6 module load seqtk
7 seqtk subseq /data/scratch/btx654/missing_busco_osedax/busco_dow
nloads/lineages/metazoa_odb10/ancestral missing_busco_list.txt
> missing_busco.fa
```


8

```
1 module load anaconda3
2 source activate augustus
3
4 gffread -w osedax_mRNA.fa -g /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax/annotation/softmasking/osedax_softmasked.fa /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax/annotation/New_annotation_Dec2020/step6/osedax_annotation_v101220.gff3
5 gffread -y osedax_protein.fa -g /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax/annotation/softmasking/osedax_softmasked.fa /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax/annotation/New_annotation_Dec2020/step6/osedax_annotation_v101220.gff3
```

blastp.sh

```
1 #!/bin/bash
2 #$ -cwd
3 #$ -j y
4 #$ -pe smp 8
5 #$ -l h_vmem=10G
6 #$ -l h_rt=120:0:0
7 #$ -l highmem
8
9 module load blast+
10 makeblastdb -in ../osedax_protein.fa -dbtype prot -out osedax_prot
11 blastp -db osedax_prot -query ../missing_busco.fa -out osedax_prot_blastp_out -max_target_seqs 5 -evalue 1e-10 -num_threads 8 -o utfmt 6
12 blastp -db osedax_prot -query ../missing_busco.fa -out osedax_prot_blastp_out.html -max_target_seqs 5 -evalue 1e-10 -num_threads
```

```
8 -html
```

blastp.sh

```
1 #!/bin/bash
2 #$ -cwd
3 #$ -j y
4 #$ -pe smp 8
5 #$ -l h_vmem=10G
6 #$ -l h_rt=120:0:0
7 #$ -l highmem
8
9 module load blast+
10 makeblastdb -in ../osedax_annotation.prot.fa -dbtype prot -out o
   sedax_prot
11 blastp -db osedax_prot -query ../missing_busco.fa -out osedax_pr
   ot_blastp_out -max_target_seqs 5 -evalue 1e-10 -num_threads 8 -o
   utfmt 6
12 blastp -db osedax_prot -query ../missing_busco.fa -out osedax_pr
   ot_blastp_out.html -max_target_seqs 5 -evalue 1e-10 -num_threads
   8 -html
```

panther.sh

```
1 #!/bin/bash
2 #$ -cwd
3 #$ -j y
4 #$ -pe smp 10
5 #$ -l h_vmem=5G
6 #$ -l h_rt=80:00:0
7 #$ -l highmem
8
```

```
9 module load perl
10 module load hmmer/
11
12 export PERL5LIB=/data/SBCS-MartinDuranLab/03-Giacomo/src/hmm scoring/lib/
13
14 perl /data/SBCS-MartinDuranLab/03-Giacomo/src/hmm scoring/panther
Score2.2.pl -l /data/SBCS-MartinDuranLab/03-Giacomo/src/hmm scoring/PANTHER15.0/ -D B -n -o panther_output -i ../missing_busco.fa
-c 10 -V -s
```

There are 23 missing Busco Panther IDs matching with annotations in Osedax: only 20 are uniq

Let's try to search for the missing proteins in Oasisia and use its proteins to blast against Osedax:

```
1 #53 different Panther IDs found in Oasisia belonging to the miss
ing Busco in osedax
2 grep -f panther_ids_list /data/SBCS-MartinDuranLab/03-Giacomo/da
ta/oasisia/annotation/New_annotation_Dec2020/step9/oasisia_annot
ation_Jan2021_TrinoPantherK0.xls | cut -f 19 | sort | uniq | wc
-l
3 cut -f 18,19 /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/a
nnotation/New_annotation_Dec2020/step9/oasisia_annotation_Jan202
1_TrinoPantherK0.xls | grep -f panther_ids_list | sort -u -k2,2
| cut -f 1 > oasisia_geneIDs_missingBUSCO_osedax.txt
4
5 module load seqtk
6 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/O
alv.fa oasisia_geneIDs_missingBUSCO_osedax_0K.txt > oasisia_prot
eins_missingBUSCO_osedax.fa
```

blastp.sh

```

1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -pe smp 8
5  #$ -l h_vmem=10G
6  #$ -l h_rt=120:0:0
7  #$ -l highmem
8
9  module load blast+
10
11 blastp -db ../../BLAST/osedax_prot -query oasisia_proteins_missingBUSCO_osedax.fa -out Oalv_VS_Ofra_missingBUSCO_blastp_out -max_target_seqs 5 -evaluate 1e-10 -num_threads 8 -outfmt 6
12 blastp -db ../../BLAST/osedax_prot -query oasisia_proteins_missingBUSCO_osedax.fa -out Oalv_VS_Ofra_missingBUSCO_blastp_out.html -max_target_seqs 5 -evaluate 1e-10 -num_threads 8 -html

```

- 22 Oasisia proteins are matching with Osedax

RECAP

62 missing BUSCO in Osedax annotation

954 tot BUSCO metazoadb10 (892 found in osedax)

BUSCO C (91.6%) + BUSCO F (1.9%) = 93.5%

METHOD	NUMBER MATCHING OSEDAX PROTEINS	percentage of missing BUSCO found
BLAST missing BUSCO (63 sequences from ancestral in metazoa_db10)	22	2.3%
BLAST missing BUSCO (53 sequences from Oasisia, obtained with PantherIDs of the missing	22	2.3%

BUSCO)		
PANTHER missing BUSCO searching directly in the annotations of Osedax	23 (actually there are just 20 unique ones)	2.4% (2.1%)
summing the three methods (see below for info)	26	2.7%

If I check the number of unique genes identified in the previous three methods:

```

1 cat BLAST/geneIDs_found_BLASTmetazodb10 Panther/geneIDs_found_P
  anther Panther/oasisia_blast/geneIDs_found_BLASToasisia | sort |
  uniq | wc -l
2 #52 but there was probably an error with panther IDs (had to cro
  p the last part of two of them in oasisia to match the ones in b
  usco missing genes)
3 cat BLAST/BUSCO_IDs_blast_metazodb10_osedax Panther/BUSCO_IDs_p
  anther_osedax Panther/oasisia_blast/BUSCO_IDs_blast_oasisia_osed
  ax | sort -u | wc -l
4 #54 different Buscos have a match in osedax
  • 52 hits, 3 of these are isoforms: 49 different genes

```

I have checked the output of busco and there are not repeated genes for different matches with busco genes. so is actually quite accurate to look at the different buscos

when i check the top blast match oasisia_vs_osedax I end up having 22 uniq osedax proteins

```

1 sort -u -k1,1 Oalv_VS_Ofra_missingBUSCO_blastp_out | cut -f 2 |
  sort -u > uniq_geneIDs_BLAST_oasisia #22 unique proteins of osed
  ax
2
3 sort -u -k1,1 osedax_prot_blastp_out | cut -f 2 | sort -u > uniq
  _geneIDs_BLAST_metazodb10 #22 unique proteins of osedax
4 sort -u -k1,1 osedax_prot_blastp_out | cut -f 1,2 > uniq_geneIDs
  _BLAST_metazodb10_vs_BUSCO

```

```
5
6 cut -f 18,19 23_missingBUSCO_found_osedax.xls | grep -f panther_
ids_list_osedax | sort -k2,2 | cut -f 1 > 23_geneIDs_osedax_Pant
her
7 cut -f 18,19 23_missingBUSCO_found_osedax.xls | grep -f panther_
ids_list_osedax | sort -k2,2 | cut -f 2 > 23_pantherIDs_osedax_P
anther
8 nano 23_geneIDs_osedax_Panther_edited > uniq_geneIDs_panther_vs_
BUSCO
9
10 grep -f panther_IDs_BUSCO_oasisia ../panther_output | cut -f1,2
| sort -k2,2 > 53_BUSCO_vs_panther_oasisia
11
```

```
1 awk -F"|" 'NR==FNR{a[$1]=$0;next} ($1 in a){b=$1;$1="";print a
[b] $0}' OFS="|" file1 file2
2 awk 'NR==FNR{a[$2]=$0;next} ($1 in a){b=$2;$1="";print a[b]
$0}' 62_metazoadb10_vs_panther 63_oasisia_geneIDs_vs_panther
```

```
1
2 join -1 2 -2 2 62_metazoadb10_vs_panther 63_oasisia_geneIDs_vs_p
anther
```

```
1 while read line; do
2 echo $line
3 annotations=$(cut -f 1,2 63_oasisia_geneIDs_vs_panther | fgrep
$line)
4 cat $annotations
5 kallisto_body=$(cut -f 2 <<< $annotations)
6 cat $kallisto_body
7 kallisto_roots=$(cut -f 3 <<< $annotations)
8 cat $kallisto_roots
9 echo $kallisto_body$'\t'$kallisto_roots >> oasisia_geneIDs_missi
```

```
ngBUSCO_found_osedax_vs_panther
```

```
10 done < oasisia_geneIDs_missingBUSCO_found_osedax
```

```
1 cut -f 1 oasisia_geneIDs_missingBUSCO_found_osedax_vs_panther |  
sed -e 's/ /\t/g' | sort -k2,2 > oasisia_geneIDs_missingBUSCO_fo  
und_osedax_vs_panther_OK  
2 cut -f 2 oasisia_geneIDs_missingBUSCO_found_osedax_vs_panther_OK  
> oasisia_geneIDs_missingBUSCO_found_osedax_vs_panther_OK_panthe  
rIDs  
3 grep -f oasisia_geneIDs_missingBUSCO_found_osedax_vs_panther_OK_  
pantherIDs 62_metazoadb10_vs_panther | sort -k2,2 > missingBUSCO  
_found_osedax_vs_panther_OK  
4  
5 paste missingBUSCO_found_osedax_vs_panther_OK oasisia_geneIDs_mi  
ssingBUSCO_found_osedax_vs_panther_OK > link_oasisia_geneIDs_BUS  
CO  
6  
7 cut -f 1 0alv_VS_Ofra_missingBUSCO_blastp_out | sed 's/0alv_/'  
> blast_output_firstColumn
```

```
1 while read line; do  
2 echo $line  
3 annotations=$(cut -f 1,3 link_oasisia_geneIDs_BUSCO | fgrep $lin  
e | cut -f 1)  
4 cat $annotations  
5 echo $annotations >> blast_output_firstColumn_vs_metazoadb10  
6 done < blast_output_firstColumn
```

```
1 cut -f 2 0alv_VS_Ofra_missingBUSCO_blastp_out > blast_output_fir  
stColumn_osedax_geneIDs  
2 paste blast_output_firstColumn_vs_metazoadb10 blast_output_first  
Column_osedax_geneIDs > uniq_geneIDs_BLAST_oasisia_vs_BUSCO  
3
```

```
4 cat BLAST/uniq_geneIDs_BLAST_metazoadb10_vs_BUSCO Panther/uniq_geneIDs_panther_vs_BUSCO Panther/oasisia_blast/uniq_geneIDs_BLAST_oasisia_vs_BUSCO > ALL_matches_osedax_vs_metazoadb10
5
6 sort -u -k1,1 ALL_matches_osedax_vs_metazoadb10 | sort -u -k1,2 | wc -l #26
```

- 26 BUSCO found in osedax (double checked for duplicates and stuff, so this is final!)

```
1 while read line; do
2   echo $line
3   annotations=$(cut -f 1,2 uniq_geneIDs_BLAST_oasisia_vs_BUSCO | f
4   grep $line | head -1)
5   cat $annotations
6   echo $annotations >> blast_output_firstColumn_vs_metazoadb10
7 done < blast_output_firstColumn
```

```
1 while read line; do
2   echo $line
3   echo "blast"
4   blast=$(cut -f 1,2 uniq_geneIDs_BLAST_metazoadb10_vs_BUSCO | fg
5   rep $line)
6   echo $blast
7   echo "blast_oasisia"
8   blast_oasisia=$(cut -f 1,2 uniq_geneIDs_BLAST_oasisia_vs_BUSCO |
9   fgrep $line)
10  echo $blast_oasisia
11  echo "panther"
12  panther=$(cut -f 1,2 uniq_geneIDs_panther_vs_BUSCO | fgrep $lin
13  e)
14  echo $panther
```



```
12 echo
   "-----
   -----
   --"
13 done < 26_BUSCO_IDs_found_osedax
```

```
1 while read line; do
2   echo $line
3   echo "Panther"
4   Panther=$(cut -f 1,2,3 ../Panther/panther_output | fgrep $line)
5   echo $Panther
6   echo
   "-----
   -----
   --"
7 done < 62_BUSCO_IDs
```

Comparison with Riftia genomes

- ✓ 1 ~~Table our Riftia vs the other Riftia (genome stats)~~
- ✓ 2 ~~Assembly vs Assembly using Minimap2 and plot~~
- ✓ 3 ~~Annotation vs Annotation BBH~~
- ✓ 4 ~~our annotation vs previous other transcriptome 2019 BBH~~
- ✓ 5 ~~PFAM barplot our annotations, riftia ann and previous transcriptome~~

1 - Table our Riftia vs the other Riftia (genome stats)

i am using:

```
1 OUR riftia_softmasked.fa
2 THEIR 4.2_RIFPA_polished_softMasked_purged_genome_v1.fasta
```

```
3 ANNOTATION_GFF3 RIFPA_final_gene_models_AUGUSTUS_v1.gff3
4 TOT_PROTEOME RIFPA_final_gene_models_AUGUSTUS_v1.prot.fasta
```

stats.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/riftia_oliveira_2022
3 #$ -o /data/scratch/btx654/riftia_oliveira_2022
4 #$ -j y
5 #$ -pe smp 4
6 #$ -l h_vmem=5G
7 #$ -l h_rt=24:0:0
8 #$ -l highmem
9
10 module load anaconda3
11 conda activate agat_env
12
13 agat_convert_sp_gxf2gxf.pl --gff RIFPA_final_gene_models_AUGUSTU
14 S_v1.gff3 --merge_loci -o RIFPA_final_gene_models_AUGUSTUS_v1_lo
15 ciMerged.gff
16
17 agat_sp_keep_longest_isoform.pl --gff RIFPA_final_gene_models_AU
18 GUSTUS_v1_lociMerged.gff -o RIFPA_final_gene_models_AUGUSTUS_v1_
19 lociMerged_longestIsoform.gff
20
21 agat_convert_sp_gxf2gxf.pl -g $final_pasa_gtf -o $final_pasa_gff
22 3
23
24 agat_sq_stat_basic.pl -i RIFPA_final_gene_models_AUGUSTUS_v1.gff
25 3 -g 4.2_RIFPA_polished_softMasked_purged_genome_v1.fasta
26
27 conda deactivate
28 source activate augustus
```

```
21
22 gffread -E RIFPA_final_gene_models_AUGUSTUS_v1_lociMerged_longes
tIsoform.gff -g 4.2_RIFPA_polished_softMasked_purged_genome_v1.f
asta -y RIFPA_final_gene_models_AUGUSTUS_v1.prot.longestIsoform.
fasta
23
24 conda deactivate
25
26 source activate quast
27
28
29 quast ./4.2_RIFPA_polished_softMasked_purged_genome_v1.fasta -o
./quast_4.2_softmasked --eukaryote
30
31
32 conda deactivate
33 source activate busco_env
34 #export BUSCO_CONFIG_FILE="/data/home/btx654/.conda/envs/busco_e
nv/busco/config/myconfig.ini"
35 #export AUGUSTUS_CONFIG_PATH=/data/SBCS-MartinDuranLab/02-Chema/
src/Augustus/config/
36
37 #busco -i RIFPA_final_gene_models_AUGUSTUS_v1.prot.longestIsofor
m.fasta -m proteins -o busco_longest_isoform -c 4 -l metazoa_odb
10
38 busco -i RIFPA_final_gene_models_AUGUSTUS_v1.prot.fasta -m prote
ins -o busco_gene_models -c 4 -l metazoa_odb10
39 busco -i 4.2_RIFPA_polished_softMasked_purged_genome_v1.fasta -m
DNA -o busco_genome -c 4 -l metazoa_odb10

1 #!/bin/bash
```

```
2  #$ -wd /data/scratch/btx654/riftia_oliveira_2022
3  #$ -o /data/scratch/btx654/riftia_oliveira_2022
4  #$ -j y
5  #$ -pe smp 4
6  #$ -l h_vmem=5G
7  #$ -l h_rt=24:0:0
8  #$ -l highmem
9
10
11 module load anaconda3
12 source activate augustus
13
14 gffread -E RIFPA_final_gene_models_AUGUSTUS_v1.gff3 -g 4.2_RIFPA
   _polished_softMasked_purged_genome_v1.fasta -y RIFPA_final_gene_
   models_AUGUSTUS_v1_protein_giacomoGFFread.fasta
15
16 conda deactivate
17 source activate busco_env
18
19 busco -i RIFPA_final_gene_models_AUGUSTUS_v1_protein_giacomoGFFr
   ead.fasta -m proteins -o busco_gene_models_giacomoGFFread -c 4 -
   l metazoa_odb10
```

	<i>OUR R. pachyptila.</i>	<i>OLIVEIRA 22</i>
Genome Size (Mb)	554	560
Number of contigs	918	447
Contig N50 (Kb)	1,424	2,870
GC content (%)	41.05	40.94
Repeats (%)	27.87	29.9

Number of genes	37,037	25977 (they say 25,984)
Number of transcripts	38,179	58020
Mean gene size (bp)	8311.47	15432.63
Mean transcript size (bp) [CM1] [CM2]	8889.27	15779.25
Gene density (per Mb)	66.85	46.39
N's	4,071	0
Busco assembly (%)	95.6	96.7
Busco annotation (%)	96.8	97.7

	Complete	Single	Duplicated	Fragmented	Missing
<i>R. pachyptila</i> assembly	95.6%	94.5%	1.1%	1%	3.4%
<i>R. pachyptila</i> annotation	96.8%	96.4%	0.4%	1%	2.2%
OLIVEIRA 22 assembly	96.7%	96.1%	0.6%	1.5%	1.8%
OLIVEIRA 22 annotation	97.7%	42.7%	55.0%	1.7%	0.6%
RIFPA_final_gene_models_AUGUSTUS_v1.prot.fasta					

```
awk '/^>/ {printf("\n%s\n",$0);next; } { printf("%s",$0);} END
{printf("\n");}' < 4.2_RIFPA_polished_softMasked_purged_genome_v
1.fasta | tail -n +2 > 4.2_RIFPA_polished_softMasked_purged_geno
me_v1_SINGLE.fasta
```

stats.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/riftia_oliveira_2022_SINGLE
3  #$ -o /data/scratch/btx654/riftia_oliveira_2022_SINGLE
4  #$ -j y
5  #$ -pe smp 4
6  #$ -l h_vmem=5G
7  #$ -l h_rt=24:0:0
8  #$ -l highmem
9
10 module load anaconda3
11 conda activate agat_env
12
13
14
15 agat_sq_stat_basic.pl -i RIFPA_final_gene_models_AUGUSTUS_v1.gff
3 -g 4.2_RIFPA_polished_softMasked_purged_genome_v1_SINGLE.fasta
16
17 conda deactivate
18 source activate augustus
19
20
21 gffread -y RIFPA_final_gene_models_AUGUSTUS_v1.prot.longestIsofo
rm.fasta -g 4.2_RIFPA_polished_softMasked_purged_genome_v1_SINGL
E.fasta RIFPA_final_gene_models_AUGUSTUS_v1_lociMerged_longestIs
```

```
oform.gff
22 gffread RIFPA_final_gene_models_AUGUSTUS_v1.gff3 -y RIFPA_final_
   gene_models_AUGUSTUS_v1_protein_giacomoGFFread.fasta -g 4.2_RIFP
   A_polished_softMasked_purged_genome_v1_SINGLE.fasta
23
24 conda deactivate
25 source activate busco_env
26 #export BUSCO_CONFIG_FILE="/data/home/btx654/.conda/envs/busco_e
   nv/busco/config/myconfig.ini"
27 #export AUGUSTUS_CONFIG_PATH=/data/SBCS-MartinDuranLab/02-Chema/
   src/Augustus/config/
28
29 #busco -i RIFPA_final_gene_models_AUGUSTUS_v1.prot.longestIsofor
   m.fasta -m proteins -o busco_longest_isoform -c 4 -l metazoa_odb
   10
30 busco -i RIFPA_final_gene_models_AUGUSTUS_v1.prot.fasta -m prote
   ins -o busco_gene_models -c 4 -l metazoa_odb10
31 busco -i 4.2_RIFPA_polished_softMasked_purged_genome_v1_SINGLE.f
   asta -m genome -o busco_genome -c 4 -l metazoa_odb10
32 busco -i RIFPA_final_gene_models_AUGUSTUS_v1_protein_giacomoGFFr
   ead.fasta -m proteins -o busco_gene_models_giacomoGFFread -c 4 -
   l metazoa_odb10
```

stats.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/riftia_oliveira_2022
3 #$ -o /data/scratch/btx654/riftia_oliveira_2022
4 #$ -j y
5 #$ -pe smp 4
6 #$ -l h_vmem=5G
7 #$ -l h_rt=24:0:0
```

```
8  #$ -l highmem
9
10 module load anaconda3
11 conda activate agat_env
12
13 agat_convert_sp_gxf2gxf.pl --gff RIFPA_final_gene_models_AUGUSTU
  S_v1.gff3 --merge_loci -o RIFPA_final_gene_models_AUGUSTUS_v1_AG
  AT.gff3
14 agat_convert_sp_gxf2gxf.pl --gff RIFPA_final_gene_models_AUGUSTU
  S_v1_AGAT.gff3 --merge_loci -o RIFPA_final_gene_models_AUGUSTUS_
  v1_AGAT_lociMerged.gff
15 agat_sp_keep_longest_isoform.pl --gff RIFPA_final_gene_models_AU
  GUSTUS_v1_AGAT_lociMerged.gff -o RIFPA_final_gene_models_AUGUSTU
  S_v1_AGAT_lociMerged_longestIsoform.gff
16
17 agat_sq_stat_basic.pl -i RIFPA_final_gene_models_AUGUSTUS_v1_AGA
  T.gff3 -g 4.2_RIFPA_polished_softMasked_purged_genome_v1.fasta
18
19 conda deactivate
20 source activate augustus
21
22 gffread -E RIFPA_final_gene_models_AUGUSTUS_v1_lociMerged_longes
  tIsoform.gff -g 4.2_RIFPA_polished_softMasked_purged_genome_v1.f
  asta -y RIFPA_final_gene_models_AUGUSTUS_v1.prot.longestIsoform.
  fasta
23
24 conda deactivate
25
26 source activate quast
27
28
29 quast ./4.2_RIFPA_polished_softMasked_purged_genome_v1.fasta -o
```



```
./quast_4.2_softmasked --eukaryote
```

```
conda deactivate
```

```
source activate busco_env
```

```
#export BUSCO_CONFIG_FILE="/data/home/btx654/.conda/envs/busco_env/busco/config/myconfig.ini"
```

```
#export AUGUSTUS_CONFIG_PATH="/data/SBCS-MartinDuranLab/02-Chema/src/Augustus/config/"
```

```
#busco -i RIFPA_final_gene_models_AUGUSTUS_v1.prot.longestIsoform.fasta -m proteins -o busco_longest_isoform -c 4 -l metazoa_odb10
```

```
busco -i RIFPA_final_gene_models_AUGUSTUS_v1.prot.fasta -m proteins -o busco_gene_models -c 4 -l metazoa_odb10
```

```
busco -i 4.2_RIFPA_polished_softMasked_purged_genome_v1.fasta -m DNA -o busco_genome -c 4 -l metazoa_odb10
```

gffread.sh

```
#!/bin/bash
```

```
#$ -wd /data/scratch/btx654/riftia_oliveira_2022_SINGLE
```

```
#$ -o /data/scratch/btx654/riftia_oliveira_2022_SINGLE
```

```
#$ -j y
```

```
#$ -pe smp 4
```

```
#$ -l h_vmem=5G
```

```
#$ -l h_rt=24:0:0
```

```
#$ -l highmem
```

```
module load anaconda3
```

```
source activate augustus
```

```
14 gffread -y RIFPA_final_gene_models_AUGUSTUS_v1_AGAT.prot.longest
    Isoform.fasta -g 4.2_RIFPA_polished_softMasked_purged_genome_v1.
    fasta RIFPA_final_gene_models_AUGUSTUS_v1_AGAT_lociMerged_longes
    tIsoform.gff
```

1	Type (3rd column)	Number	Size total (kb)	Si
	ze mean (bp)	% of the genome	/!\Results are roundi	ng to two decimal places
2	cds	496239	103636.29	208.84
3	exon	574346	103714.43	180.58
4	five_prime_utr	35671	35.67	1.00
5	gene	25977	400893.47	15432.63
6	mrna	58020	915511.98	15779.25
7	start_codon	57895	173.69	3.00
8	stop_codon	57926	173.78	3.00
9	three_prime_utr	42439	42.47	1.00
10	tss	35670	35.67	1.00
11	tts	42424	42.42	1.00
12	Total	1426607	1524259.86	1068.45
				27
				1.81

2 - Assembly vs Assembly using Minimap2 and plot

asm5/asm10/asm20: asm-to-ref mapping, for ~0.1/1/5% sequence divergence. asm5 for the same species

```
/minimap2 -cx asm5 asm1.fa asm2.fa > aln.paf # intra
-species asm-to-asm alignment
```

```
minimap2 -ax sr $pacbio_corrected_nonBacteria $R1_cleaned $R2_cleaned --split-prefix temp_sam_ > $alignment_sam
```

```
minimap2 -x map-pb $ref_genome $pb_fasta | gzip -c - > purge_${1}.paf.gz
```

QUERY ours
TARGET theirs

minimap2.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/alignment
3  #$ -o /data/scratch/btx654/alignment
4  #$ -j y
5  #$ -pe smp 12
6  #$ -l h_vmem=30G
7  #$ -l h_rt=240:0:0
8  #$ -l highmem
9
10 module load anaconda3
11 source activate minimap2_env
12
13 minimap2 -cx asm5 /data/scratch/btx654/riftia_oliveira_2022_SINGLE/4.2_RIFPA_polished_softMasked_purged_genome_v1_SINGLE.fasta /data/SBCS-MartinDuranLab/03-Giacomo/data/riftia/annotation/riftia_softmasked.fa > alignment_riftia_OUR_VS_oliveira_2022.paf
```

```
scp -i ~/.ssh/id_rsa_apocrita -r btx654@login.hpc.qmul.ac.uk:/data/scratch/btx654/alignment/alignment_riftia_OUR_VS_oliveira_2022.paf /Users/giacomo/Desktop/
```

R

```
install.packages("pafr")
```

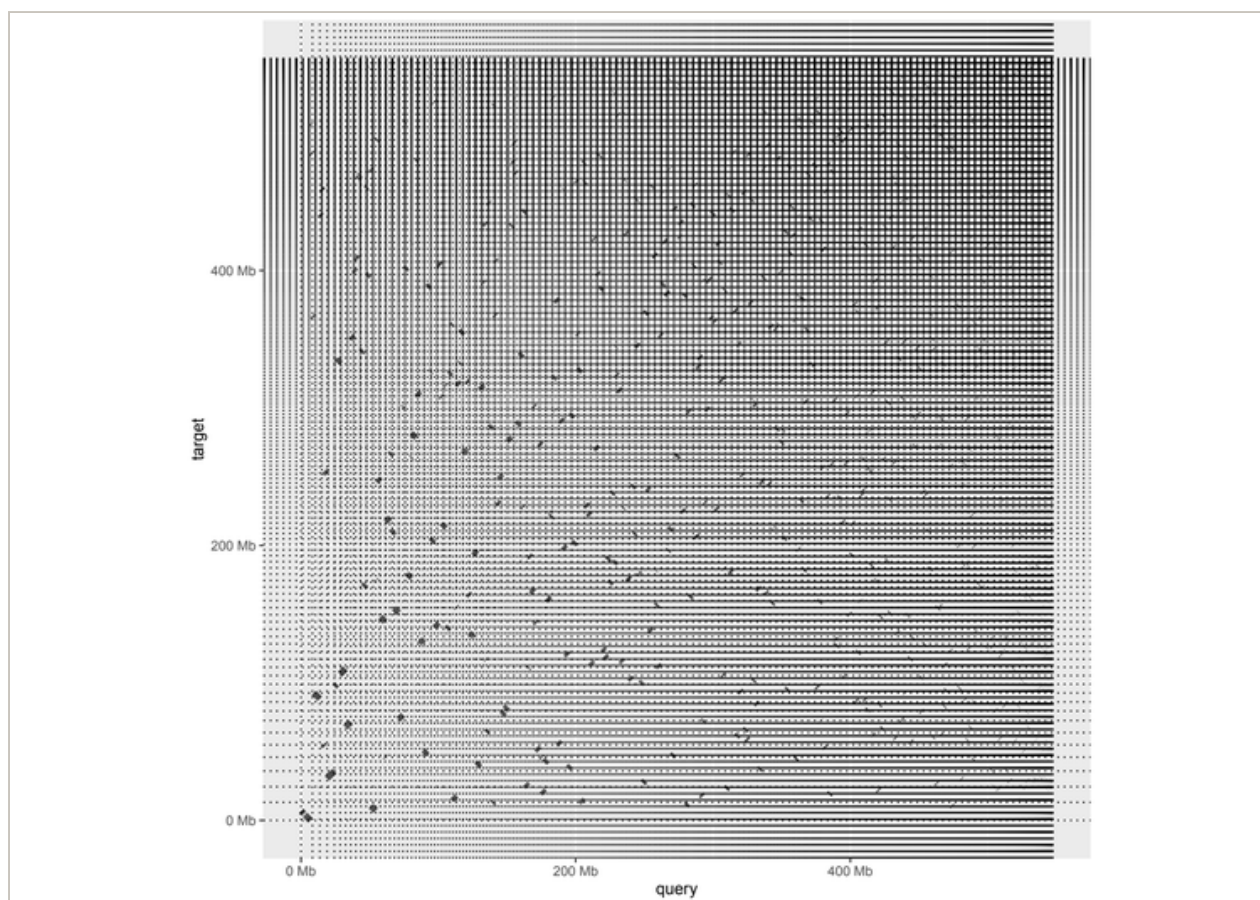
```
1 library(pafr)
2 library(ggplot2)
3 ## Loading required package: ggplot2
4 path_to_alignment <- system.file("extdata", "/Users/giacomo/Desktop/alignment_riftia_OUR_VS_oliveira2022/alignment_riftia_OUR_VS_oliveira_2022.paf", package = "pafr")
5 ali <- read_paf(path_to_alignment)
6 dotplot(ali)
7
8 ali <- read_paf("/Users/giacomo/Desktop/alignment_riftia_OUR_VS_oliveira2022/alignment_riftia_OUR_VS_oliveira_2022.paf")
9 prim_alignment <- filter_secondary_alignments(ali)
10 long_ali <- subset(prim_alignment, alen > 1e4 & mapq > 40)
```

```
over_N50_ali <- subset(long_ali, qlen > 1423584, tlen > 2870320)
```

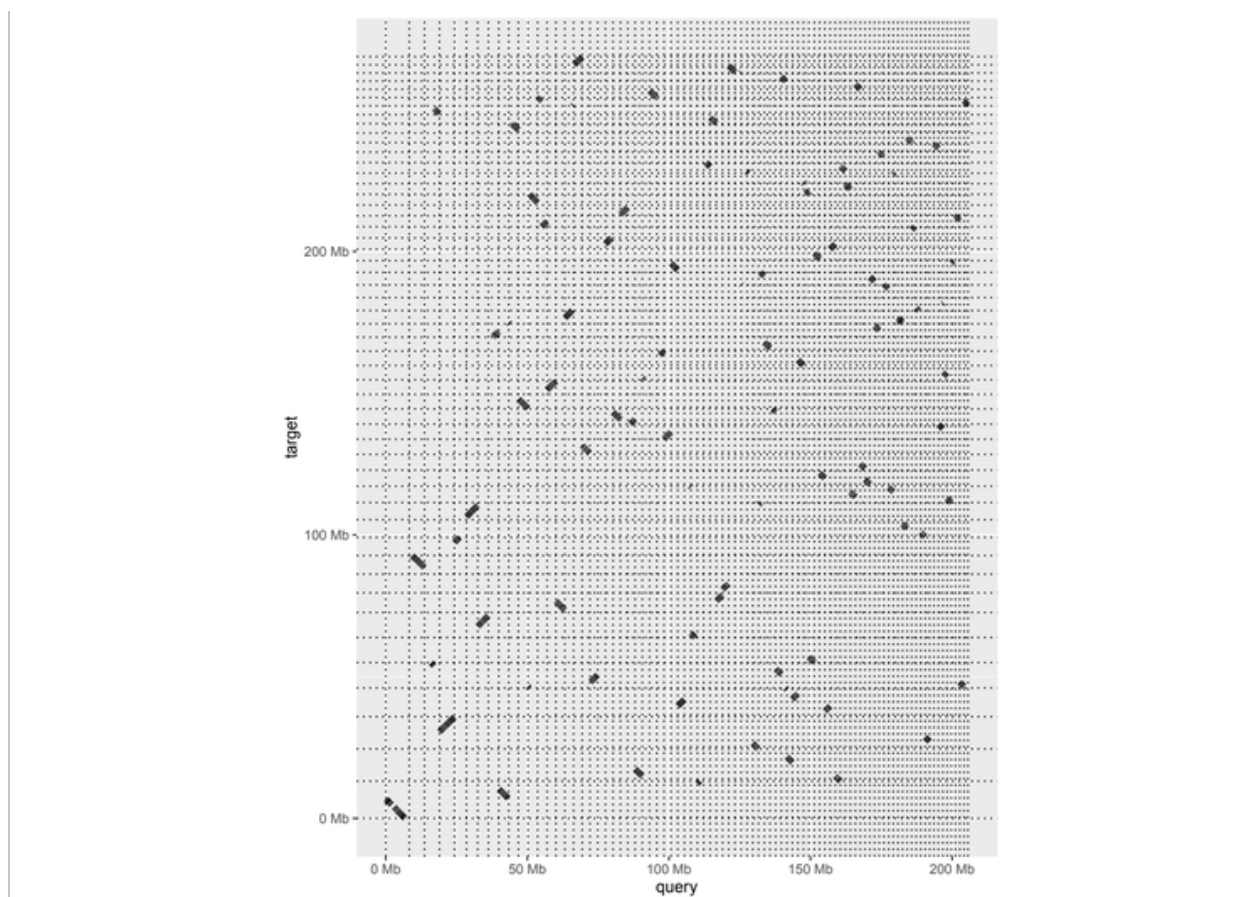
selecting just the sequences equal or over N50 size

```
awk -F "\t" '{ if(($2 >= 1423584) && ($7 >= 2870320)) { print } }' alignment_riftia_OUR_VS_oliveira_2022.paf > overN50_alignment_riftia_OUR_VS_oliveira_2022.paf
```

```
1 overN50_ali <- read_paf("/Users/giacomo/Desktop/alignment_riftia_OUR_VS_oliveira2022/overN50_alignment_riftia_OUR_VS_oliveira_2022.paf")
2 overN50_prim_alignment <- filter_secondary_alignments(overN50_ali)
3 overN50_long_ali <- subset(overN50_prim_alignment, alen > 1e4 & mapq > 40)
4 dotplot(overN50_long_ali)
```



PDF dotplot_prim_alignment_long_ali • PDF document



PDF dotplot_overN50_prim_alignment_long_ali • PDF document

3 - Annotation vs Annotation

INPUT FILES:

OUR **Rpac.fa** (Non redundant proteome, longest isoform)

OLIVEIRA2022

RIFPA_final_gene_models_AUGUSTUS_v1_AGAT.prot.longestIsoform.fasta (longest isoform generated in 1 with agat and gffread)

multi-line FASTA (default from NCBI) to single-line FASTA

```
1 awk '/^>/ {printf("\n%s\n",$0);next; } { printf("%s",$0);} END {printf("\n");}' < Rpac.fa | tail -n +2 > Rpac_SINGLE.fa
2 awk '/^>/ {printf("\n%s\n",$0);next; } { printf("%s",$0);} END {printf("\n");}' < RIFPA_final_gene_models_AUGUSTUS_v1_AGAT.pro
t.longestIsoform.fasta | tail -n +2 > oliveira2022_SINGLE.fa
```

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 8
6  #$ -l h_vmem=1G
7  #$ -l h_rt=36:0:0
8
9  species=$1
10
11  echo "Working on "$species
12
13  cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
    ec2020/step8/
14
15  module load anaconda3
16  conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
    3/trinotate_env
17
18  blastp -query input_trinotate_proteins.fa -db /data/SBCS-MartinD
    uranLab/03-Giacomo/db/trinotate/uniprot_sprot.pep -num_threads 8
    -max_target_seqs 1 -outfmt 6 -evalue 1e-3 > blastp.outfmt6
```

GUIDE

runBLAST.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/ann_VS_ann_oliveira2022/
```

```
3  #$ -o /data/scratch/btx654/ann_VS_ann_oliveira2022/
4  #$ -j y
5  #$ -pe smp 8
6  #$ -l h_vmem=10G
7  #$ -l h_rt=36:0:0
8  #$ -l highmem
9
10 #Script to perform a reciprocal blast search
11 #Usage: bash runRBLAST.sh PATH/TO/QUERY/FILE PATH/TO/DB/FILE PAT
    H/TO/OUTPUTS
12 #Usage ex: bash runRBLAST.sh PATH/TO/INPUT1/species1.fasta.trans
    decoder.pep PATH/TO/INPUT2/species2.fasta.transdecoder.pep PATH/
    TO/OUTPUTS
13
14
15 #Input query file
16 inputQuery=$1
17 #Input DB reciprocal file
18 inputDB=$2
19 #Path to output results
20 outputPath=$3
21 #Move to DB directory
22 queryPath=$(dirname $inputQuery)
23
24 module load anaconda3
25 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
    3/trinotate_env
26
27 #cd $queryPath
28 #Make blastable protein DB of the input query
29 makeblastdb -in $inputQuery -dbtype prot
```



```
30 #Move to query directory
31 #dbPath=$(dirname $inputDB)
32 #cd $dbPath
33 #Make blastable protein DB of the input DB
34 makeblastdb -in $inputDB -dbtype prot
35 #Output start status message
36 echo "Beginning reciprocal BLAST..."
37 #Move to outputs folder
38 #cd $outputPath
39 #Use blastp to search a database
40 blastp -query $inputQuery -db $inputDB -max_target_seqs 1 -outfmt
t 6 -evaluate 1e-3 -num_threads 8 > blast.outfmt6
41 #Switch query and search paths for reciprocal search
42 blastp -query $inputDB -db $inputQuery -max_target_seqs 1 -outfmt
t 6 -evaluate 1e-3 -num_threads 8 > blast_reciprocal.outfmt6
43 #Output end status message
44 echo "Finished reciprocal BLAST!"
```

```
qsub runRBLAST.sh /data/scratch/btx654/ann_VS_ann_oliveira2022/o
liveira2022_SINGLE.fa /data/scratch/btx654/ann_VS_ann_oliveira20
22/Rpac_SINGLE.fa /data/scratch/btx654/ann_VS_ann_oliveira2022/
```

findRBH.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/ann_VS_ann_oliveira2022/
3 #$ -o /data/scratch/btx654/ann_VS_ann_oliveira2022/
4 #$ -j y
5 #$ -pe smp 1
6 #$ -l h_vmem=10G
7 #$ -l h_rt=36:0:0
8 #$ -l highmem
```

```
9
10 #Script to filter reciprocal blast results for best hits
11 #Usage: bash findRBH.sh PATH/TO/QUERY/BLAST/RESULTS PATH/TO/DB/BLAST/RESULTS
12 #Usage ex: bash findRBH.sh blast.outfmt6 blast_reciprocal.outfmt6
13 #Input query blast results file
14 queryPath=$1
15 #Input DB reciprocal blast results file
16 dbPath=$2
17 #Final output files
18 outFileRBH="blast_RBH.txt"
19 outFileSummary="blast_RBH_summary.txt"
20 #Add headers to output RBH files
21 echo "queryHit,dbHit" > $outFileRBH
22 echo "queryHits,dbHits,bestHits" > $outFileSummary
23 #Output start status message
24 echo "Recording RBH..."
25 #Loop over query blast results
26 while IFS=$'\t' read -r f1 f2 f3 f4 f5 f6 f7 f8 f9 f10 f11 f12
27 do
28 #Determine RBH to DB blast results
29 if grep -q "$f2$'\t' "$f1$'\t' $dbPath; then #RBH
30 echo "$f1,$f2" >> $outFileRBH
31 fi
32 done < $queryPath
33 #Output summary of RBH
34 queryHits=$(wc -l "$queryPath" | cut -d ' ' -f 1)
35 dbHits=$(wc -l "$dbPath" | cut -d ' ' -f 1)
36 bestHits=$((($wc -l "$outFileRBH" | cut -d ' ' -f 1)-1))
37 echo "$queryHits","$dbHits","$bestHits" >> $outFileSummary
```

```
38 #Output end status message
39 echo "Finished recording RBH!"
```

```
qsub findRBH.sh blast.outfmt6 blast_reciprocal.outfmt6
```

reciprocal hits 17988

35282 proteins of our riftia had a match with the other riftia
23441 proteins of the other riftia had a match with our riftia

4 - Annotation vs previous other transcriptome 2019

```
1 conda create -n SRAtools_env
2 conda activate SRAtools_env
3 conda install -c bioconda sra-tools
```

```
fasterq-dump SRA_list.txt -O /data/scratch/btx654/riftia_hinzke_2019
```

sra.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/riftia_hinzke_2019
3 #$ -o /data/scratch/btx654/riftia_hinzke_2019
4 #$ -j y
5 #$ -pe smp 4
6 #$ -l h_vmem=5G
7 #$ -l h_rt=24:0:0
8 #$ -l highmem
9
10 module load anaconda3
11 conda activate SRAtools_env
12
13 while read line; do
```

```
14 fastq-dump --defline-seq '@$sn[_$rn]/$ri' --split-files $line
    -O /data/scratch/btx654/riftia_hinzke_2019
15 done < SRA_list.txt
```

SRA_TOOLKIT/fastq-dump --defline-seq '@\$sn[_\$rn]/\$ri' --split-files file.sra

Installation:

```
1 module load anaconda3
2 conda create -n Trinity_env
3 source activate Trinity_env
4 conda install -c bioconda trinity
```

- Trinity-v2.8.5

```
1 ace      ace_R1  owe_ace_R1_r1__paired.fastq.gz  owe_ace_R1_r2_pa
  ired.fastq.gz
2 1cell    1cell_R1      owe_1cell_R1_r1__paired.fastq.gz
  owe_1cell_R1_r2_paired.fastq.gz
3 2cell    2cell_R1      owe_2cell_R1_r1__paired.fastq.gz
  owe_2cell_R1_r2_paired.fastq.gz
4 4cell    4cell_R1      owe_4cell_R1_r1__paired.fastq.gz
  owe_4cell_R1_r2_paired.fastq.gz
5 8cell    8cell_R1      owe_8cell_R1_r1__paired.fastq.gz owe_8cel
  l_R1_r2_paired.fastq.gz
6 3h       3h_R1      owe_3h_R1_r1__paired.fastq.gz  owe_3h_R1_r2_pai
  red.fastq.gz
7 4h       4h_R1      owe_4h_R1_r1__paired.fastq.gz  owe_4h_R1_r2_pai
  red.fastq.gz
8 5h       5h_R1      owe_5h_R1_r1__paired.fastq.gz  owe_5h_R1_r2_pai
  red.fastq.gz
9 9h       9h_R1      owe_9h_R1_r1__paired.fastq.gz  owe_9h_R1_r2_pai
  red.fastq.gz
10 13h      13h_R1     owe_13h_R1_r1__paired.fastq.gz  owe_13h_R1_r2_pa
```

```
ired.fastq.gz
11 18h      18h_R1  owe_18h_R1_r1__paired.fastq.gz  owe_18h_R1_r2_pa
    ired.fastq.gz
12 27h      27h_R1  owe_27h_R1_r1__paired.fastq.gz  owe_27h_R1_r2_pa
    ired.fastq.gz
```

```
1  #!/bin/bash
2  #$ -pe smp 20
3  #$ -l highmem
4  #$ -l h_vmem=10G
5  #$ -l h_rt=240:0:0
6  #$ -cwd
7  #$ -j y
8
9  module load trinity/2.4.0
10
11 Trinity --seqType fq --max_memory 200G --samples_file tissueLibr
    aries_R1.txt --SS_lib_type RF --CPU 10 --output Oxford_Illumina_
    trinity_R1
```

```
1 Trinity \
2   --seqType fq \
3   --left $sample1_r1,$sample2_r1 \
4   --right $sample1_r2,$sample2_r2 \
5   --SS_lib_type RF \
6   --max_memory 400G \
7   --CPU 20 \
8   --output $output \
9   --full_cleanup \
10  --trimmomatic
```

1	SRR8949056	SRR8949056	SRR8949056_1.fastq	SRR8949056_2.fastq
2	SRR8949057	SRR8949057	SRR8949057_1.fastq	SRR8949057_2.fastq
3	SRR8949058	SRR8949058	SRR8949058_1.fastq	SRR8949058_2.fastq
4	SRR8949059	SRR8949059	SRR8949059_1.fastq	SRR8949059_2.fastq
5	SRR8949060	SRR8949060	SRR8949060_1.fastq	SRR8949060_2.fastq
6	SRR8949061	SRR8949061	SRR8949061_1.fastq	SRR8949061_2.fastq
7	SRR8949062	SRR8949062	SRR8949062_1.fastq	SRR8949062_2.fastq
8	SRR8949063	SRR8949063	SRR8949063_1.fastq	SRR8949063_2.fastq
9	SRR8949064	SRR8949064	SRR8949064_1.fastq	SRR8949064_2.fastq
10	SRR8949065	SRR8949065	SRR8949065_1.fastq	SRR8949065_2.fastq
11	SRR8949066	SRR8949066	SRR8949066_1.fastq	SRR8949066_2.fastq
12	SRR8949067	SRR8949067	SRR8949067_1.fastq	SRR8949067_2.fastq
13	SRR8949068	SRR8949068	SRR8949068_1.fastq	SRR8949068_2.fastq
14	SRR8949069	SRR8949069	SRR8949069_1.fastq	SRR8949069_2.fastq
15	SRR8949070	SRR8949070	SRR8949070_1.fastq	SRR8949070_2.fastq
16	SRR8949071	SRR8949071	SRR8949071_1.fastq	SRR8949071_2.fastq

```

q
17 SRR8949072      SRR8949072  SRR8949072_1.fastq  SRR8949072_2.fast
q
18 SRR8949073      SRR8949073  SRR8949073_1.fastq  SRR8949073_2.fast
q
19 SRR8949074      SRR8949074  SRR8949074_1.fastq  SRR8949074_2.fast
q
20 SRR8949075      SRR8949075  SRR8949075_1.fastq  SRR8949075_2.fast
q
21 SRR8949076      SRR8949076  SRR8949076_1.fastq  SRR8949076_2.fast
q
22 SRR8949077      SRR8949077  SRR8949077_1.fastq  SRR8949077_2.fast
q

```

combine samples all together and then trinity on the combined sample

```

1 cat SRR8949056_1.fastq SRR8949057_1.fastq SRR8949058_1.fastq SRR
  8949059_1.fastq SRR8949060_1.fastq SRR8949061_1.fastq SRR8949062
  _1.fastq SRR8949063_1.fastq SRR8949064_1.fastq SRR8949065_1.fast
  q SRR8949066_1.fastq SRR8949067_1.fastq SRR8949068_1.fastq SRR89
  49069_1.fastq SRR8949070_1.fastq SRR8949071_1.fastq SRR8949072_
  1.fastq SRR8949073_1.fastq SRR8949074_1.fastq SRR8949075_1.fastq
  SRR8949076_1.fastq SRR8949077_1.fastq > combined_1.fastq
2 cat SRR8949056_2.fastq SRR8949057_2.fastq SRR8949058_2.fastq SRR
  8949059_2.fastq SRR8949060_2.fastq SRR8949061_2.fastq SRR8949062
  _2.fastq SRR8949063_2.fastq SRR8949064_2.fastq SRR8949065_2.fast
  q SRR8949066_2.fastq SRR8949067_2.fastq SRR8949068_2.fastq SRR89
  49069_2.fastq SRR8949070_2.fastq SRR8949071_2.fastq SRR8949072_
  2.fastq SRR8949073_2.fastq SRR8949074_2.fastq SRR8949075_2.fastq
  SRR8949076_2.fastq SRR8949077_2.fastq > combined_2.fastq

```

```

1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/riftia_hinzke_2019
3 #$ -j y

```

```
4  #$ -o /data/scratch/btx654/riftia_hinzke_2019
5  #$ -pe smp 20
6  #$ -l h_vmem=20G
7  #$ -l h_rt=240:0:0
8  #$ -l highmem
9
10 output=riftia_hinzke_2019_trinity
11
12 module load anaconda3
13 source activate Trinity_env
14
15 Trinity \
16   --seqType fq \
17   --left combined_1.fastq \
18   --right combined_2.fastq \
19   --SS_lib_type RF \
20   --max_memory 400G \
21   --CPU 20 \
22   --output $output \
23   --full_cleanup \
24   --trimmomatic
25
26
27 if [ -f "$output".Trinity.fasta ]
28 then
29     if [ -s "$output".Trinity.fasta ]
30     then
31         echo $output".Trinity.fasta exists and not empty"
32     else
33         echo $output".Trinity.fasta exists but empty"
```



```
34     fi
35 else
36     echo $output".Trinity.fasta not exists"
37 fi
```

trinity_single.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/riftia_hinzke_2019
3  #$ -j y
4  #$ -o /data/scratch/btx654/riftia_hinzke_2019
5  #$ -pe smp 20
6  #$ -l h_vmem=20G
7  #$ -l h_rt=240:0:0
8  #$ -l highmem
9
10 input=$1
11 R1="$input"_1.fastq
12 R2="$input"_2.fastq
13 output="$input"_trinity
14
15 module load anaconda3
16 source activate Trinity_env
17
18 Trinity \
19   --seqType fq \
20   --left $R1 \
21   --right $R2 \
22   --SS_lib_type RF \
23   --max_memory 400G \
24   --CPU 20 \
25   --output $output \
```

```
26 --full_cleanup \  
27 --trimmomatic  
28  
29  
30 if [ -f "$output".Trinity.fasta ]  
31 then  
32     if [ -s "$output".Trinity.fasta ]  
33     then  
34         echo $output".Trinity.fasta exists and not empty"  
35     else  
36         echo $output".Trinity.fasta exists but empty"  
37     fi  
38 else  
39     echo $output".Trinity.fasta not exists"  
40 fi
```

now I need to merge all the transcriptomes into a single one with cd-hit

```
1 conda create -n cdhit_env  
2 conda activate cdhit_env  
3 conda install -c bioconda cd-hit
```

- cd hit version v4.8.1

first I need to merge all the transcriptome into one with cat

```
cat SRR8949056_trinity.Trinity.fasta SRR8949057_trinity.Trinity.  
fasta SRR8949058_trinity.Trinity.fasta SRR8949059_trinity.Trinit  
y.fasta SRR8949060_trinity.Trinity.fasta SRR8949061_trinity.Trin  
ity.fasta SRR8949062_trinity.Trinity.fasta SRR8949063_trinity.Tr  
inity.fasta SRR8949064_trinity.Trinity.fasta SRR8949065_trinity.  
Trinity.fasta SRR8949066_trinity.Trinity.fasta SRR8949067_trinit  
y.Trinity.fasta SRR8949068_trinity.Trinity.fasta SRR8949069_trin  
ity.Trinity.fasta SRR8949070_trinity.Trinity.fasta SRR8949071_tr  
inity.Trinity.fasta SRR8949072_trinity.Trinity.fasta SRR8949073_  
trinity.Trinity.fasta SRR8949074_trinity.Trinity.fasta SRR894907
```

```
5_trinity.Trinity.fasta SRR8949076_trinity.Trinity.fasta SRR8949
077_trinity.Trinity.fasta > combined_trinity.Trinity.fasta
```

- there are 2284338 sequences in combined_trinity.Trinity.fasta

cdhit.sh

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -pe smp 20
5  #$ -l h_vmem=20G
6  #$ -l h_rt=240:0:0
7  #$ -l highmem
8
9  /data/SBCS-MartinDuranLab/02-Chema/src/cdhit/cd-hit-est -i combi
ned_trinity.Trinity.fasta -o cdhit_90similarity -M 380000 -T 20
```

-i = input

-o = output

-c = cut-off

-n = word size:

n=5 for thresholds 0.7 ~ 1.0

n=4 for thresholds 0.6 ~ 0.7

n=3 for thresholds 0.5 ~ 0.6

n=2 for thresholds 0.4 ~ 0.5

-M = maximum available memory

-T threads

5 - PFAM barplot

generate a list of uniq pfam from riftia ann file:

```
1  head -5 riftia_annotation_Jan2021_TrinoPantherK0.xls | cut -f 8
   | sed '/^\./d' | sed "s/\\`/\\n/g" | sed "s/\\^.*//g"
2
3  cut -f 8 riftia_annotation_Jan2021_TrinoPantherK0.xls | sed 's/^\n'
```

```
/\./g' | sed "s/\`/\n/g" | sed "s/\^.*//g" | sort | uniq | wc -l  
4 cut -f 8 riftia_annotation_Jan2021_TrinoPantherK0.xls | sed "s/\`/\n/g" | sed "s/\^.*//g" | sort | uniq | wc -l #then remove first and last lines
```

16467 proteins annotated with pfam

5435 unique pfam

HMMER_Rpac.sh

```
1 #!/bin/bash  
2 #$ -wd /data/scratch/btx654/pfam  
3 #$ -o /data/scratch/btx654/pfam  
4 #$ -j y  
5 #$ -l highmem  
6 #$ -pe smp 12  
7 #$ -l h_vmem=40G  
8 #$ -l h_rt=36:0:0  
9 #$ -l highmem  
10  
11 module load anaconda3  
12 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/trinotate_env  
13  
14 hmmscan --cpu 12 --domtblout PFAM_Rpac.out /data/SBCS-MartinDuranLab/03-Giacomo/db/trinotate/Pfam-A.hmm Rpac.fa > pfam_Rpac.log
```

HMMER_oliveira2022NR.sh

```
1 #!/bin/bash  
2 #$ -wd /data/scratch/btx654/pfam  
3 #$ -o /data/scratch/btx654/pfam  
4 #$ -j y  
5 #$ -l highmem
```

```
6  #$ -pe smp 12
7  #$ -l h_vmem=40G
8  #$ -l h_rt=36:0:0
9  #$ -l highmem
10
11 module load anaconda3
12 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
   3/trinotate_env
13
14 hmmscan --cpu 12 --domtblout PFAM_oliveira2022NR.out /data/SBCS-
   MartinDuranLab/03-Giacomo/db/trinotate/Pfam-A.hmm RIFPA_final_ge
   ne_models_AUGUSTUS_v1_AGAT.prot.longestIsoform.fasta > pfam_oliv
   eira2022NR.log
```

check how many proteins have been annotated with pfam:

```
grep -v "^#" PFAM_oliveira2022NR.out | grep -o "\sRIFPA.*t[0-9]"
| sort | uniq | wc -l
```

20805

check how many different pfam in oliveira2022

```
grep -v "^#" PFAM_oliveira2022NR.out | grep -o "\sPF....."
| sort | uniq | wc -l
```

14439 uniq pfam

I need to filter the output in a better way, there are too many overlapping pfam

```
1 conda create -n pfamScan_env
2 conda activate pfamScan_env
3 conda install -c bioconda pfam_scan
```

```
1 pfam_scan.pl -fasta <fasta_file> -dir <directory location of Pfa
  m files>
```

```
2 Useful options are:
3 -outfile <file> : output file, otherwise send to STDOUT
```

pfamScan_oliveira2022NR.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/pfam/pfamScan
3 #$ -o /data/scratch/btx654/pfam/pfamScan
4 #$ -j y
5 #$ -l highmem
6 #$ -pe smp 12
7 #$ -l h_vmem=40G
8 #$ -l h_rt=36:0:0
9 #$ -l highmem
10
11 module load anaconda3
12 conda activate pfamScan_env
13
14 pfam_scan.pl -cpu 12 -fasta ../RIFPA_final_gene_models_AUGUSTUS_
v1_AGAT.prot.longestIsoform.fasta -dir /data/SBCS-MartinDuranLab
/00-BlastDBs -outfile PFAMscan_oliveira2022NR.out
```

pfamScan_Rpac.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/pfam/pfamScan
3 #$ -o /data/scratch/btx654/pfam/pfamScan
4 #$ -j y
5 #$ -l highmem
6 #$ -pe smp 12
7 #$ -l h_vmem=40G
8 #$ -l h_rt=36:0:0
```

```
9  #$ -l highmem
10
11  module load anaconda3
12  conda activate pfamScan_env
13
14  pfam_scan.pl -cpu 12 -fasta ../Rpac.fa -dir /data/SBCS-MartinDuranLab/00-BlastDBs -outfile PFAMscan_Rpac.out
```

check how many proteins have been annotated with pfam:

```
grep -v "^#" PFAMscan_oliveira2022NR.out | grep -o "RIFPA.*t[0-9]" | sort | uniq | wc -l
```

13179

check how many different pfam in oliveira2022

```
grep -v "^#" PFAMscan_oliveira2022NR.out | grep -o "\sPF....." | sort | uniq | wc -l
```

5079 uniq pfam

CHEMA

CD-hit merge transcriptomes

input file:

```
/data/SBCS-MartinDuranLab/03-Giacomo/data/07-Reviews/riftia_hinke_2019/combined_trinity.Trinity.fasta
```

cdhit.sh

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
```

```
4  #$ -pe smp 20
5  #$ -l h_vmem=20G
6  #$ -l h_rt=240:0:0
7  #$ -l highmem
8
9  /data/SBCS-MartinDuranLab/02-Chema/src/cdhit/cd-hit-est -i combined_trinity.Trinity.fasta -o cdhit_90similarity -M 380000 -T 20
```

Best Blast hit

[GUIDE](#)

use the merged cd-hit transcriptome as input

Our riftia vs oliveira 2022 directory with all the db files generated:

```
/data/SBCS-MartinDuranLab/03-Giacomo/data/07-Reviews/riftia_oliveira_2022/ann_VS_ann_oliveira2022/
```

runBLAST.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/ann_VS_ann_oliveira2022/
3  #$ -o /data/scratch/btx654/ann_VS_ann_oliveira2022/
4  #$ -j y
5  #$ -pe smp 8
6  #$ -l h_vmem=10G
7  #$ -l h_rt=36:0:0
8  #$ -l highmem
9
10 #Script to perform a reciprocal blast search
11 #Usage: bash runRBLAST.sh PATH/TO/QUERY/FILE PATH/TO/DB/FILE PATH/TO/OUTPUTS
12 #Usage ex: bash runRBLAST.sh PATH/TO/INPUT1/species1.fasta.trans
```



```
decoder.pep PATH/T0/INPUT2/species2.fasta.transdecoder.pep PATH/
T0/OUTPUTS

13
14
15 #Input query file
16 inputQuery=$1
17 #Input DB reciprocal file
18 inputDB=$2
19 #Path to output results
20 outputPath=$3
21 #Move to DB directory
22 queryPath=$(dirname $inputQuery)
23
24 module load anaconda3
25 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
3/trinotate_env
26
27 #cd $queryPath
28 #Make blastable protein DB of the input query
29 makeblastdb -in $inputQuery -dbtype prot
30 #Move to query directory
31 #dbPath=$(dirname $inputDB)
32 #cd $dbPath
33 #Make blastable protein DB of the input DB
34 makeblastdb -in $inputDB -dbtype prot
35 #Output start status message
36 echo "Beginning reciprocal BLAST..."
37 #Move to outputs folder
38 #cd $outputPath
39 #Use blastp to search a database
40 blastp -query $inputQuery -db $inputDB -max_target_seqs 1 -outfm
```

```
t 6 -evaluate 1e-3 -num_threads 8 > blast.outfmt6
41 #Switch query and search paths for reciprocal search
42 blastp -query $inputDB -db $inputQuery -max_target_seqs 1 -outfmt
t 6 -evaluate 1e-3 -num_threads 8 > blast_reciprocal.outfmt6
43 #Output end status message
44 echo "Finished reciprocal BLAST!"
```

```
qsub runRBLAST.sh /data/scratch/btx654/ann_VS_ann_oliveira2022/o
liveira2022_SINGLE.fa /data/scratch/btx654/ann_VS_ann_oliveira20
22/Rpac_SINGLE.fa /data/scratch/btx654/ann_VS_ann_oliveira2022/
```

findRBH.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/ann_VS_ann_oliveira2022/
3 #$ -o /data/scratch/btx654/ann_VS_ann_oliveira2022/
4 #$ -j y
5 #$ -pe smp 1
6 #$ -l h_vmem=10G
7 #$ -l h_rt=36:0:0
8 #$ -l highmem
9
10 #Script to filter reciprocal blast results for best hits
11 #Usage: bash findRBH.sh PATH/TO/QUERY/BLAST/RESULTS PATH/TO/DB/B
LAST/RESULTS
12 #Usage ex: bash findRBH.sh blast.outfmt6 blast_reciprocal.outfmt
6
13 #Input query blast results file
14 queryPath=$1
15 #Input DB reciprocal blast results file
16 dbPath=$2
17 #Final output files
18 outFileRBH="blast_RBH.txt"
```

```
19 outFileSummary="blast_RBH_summary.txt"
20 #Add headers to output RBH files
21 echo "queryHit,dbHit" > $outFileRBH
22 echo "queryHits,dbHits,bestHits" > $outFileSummary
23 #Output start status message
24 echo "Recording RBH..."
25 #Loop over query blast results
26 while IFS=$'\t' read -r f1 f2 f3 f4 f5 f6 f7 f8 f9 f10 f11 f12
27 do
28 #Determine RBH to DB blast results
29 if grep -q "$f2$'\t'"$f1$'\t' $dbPath; then #RBH
30 echo "$f1,$f2" >> $outFileRBH
31 fi
32 done < $queryPath
33 #Output summary of RBH
34 queryHits=$(wc -l "$queryPath" | cut -d ' ' -f 1)
35 dbHits=$(wc -l "$dbPath" | cut -d ' ' -f 1)
36 bestHits=$((($wc -l "$outFileRBH" | cut -d ' ' -f 1)-1))
37 echo "$queryHits","$dbHits","$bestHits" >> $outFileSummary
38 #Output end status message
39 echo "Finished recording RBH!"
```

```
qsub findRBH.sh blast.outfmt6 blast_reciprocal.outfmt6
```

reciprocal hits 17988

reciprocal hits (transcriptome vs our Riftia) = 15469

I re-do this using the script: <https://scriptomika.wordpress.com/2014/01/28/extract-best-reciprocal-blast-matches/>

```
../get_RBH.py blast.outfmt6 blast_reciprocal.outfmt6 1 2 11 low
OliveiraVSGenome.hits.out
```

Genome vs mBIO == 15469

Genome vs Oliveira == 17981

35282 proteins of our riftia had a match with the other riftia

23441 proteins of the other riftia had a match with our riftia

PFAM

use the merged cd-hit transcriptome as input

pfamScan_Rpac.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/pfam/pfamScan
3  #$ -o /data/scratch/btx654/pfam/pfamScan
4  #$ -j y
5  #$ -l highmem
6  #$ -pe smp 12
7  #$ -l h_vmem=40G
8  #$ -l h_rt=36:0:0
9  #$ -l highmem
10
11 module load anaconda3
12 conda activate pfamScan_env
13
14 pfam_scan.pl -cpu 12 -fasta ../Rpac.fa -dir /data/SBCS-MartinDur
  anLab/00-BlastDBs -outfile PFAMscan_Rpac.out
```

check how many proteins have been annotated with pfam:

```
grep -v "^#" PFAMscan_oliveira2022NR.out | grep -o "RIFPA.*t[0-9]" | sort | uniq | wc -l
```

13179

Trinity CD-HIT: 5187 unique PFAM (by Chema)

chema's code:

```
grep -v "^#" Rpac_cdhit_transD_pfamscan.out | tr "[:space:]" "\n" | grep -E '^PF[0-9]' | sort | uniq | wc -l
```

Chapter 4

Gene family analyses

Broccoli

```
1 conda create --prefix /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/broccoli_env python=3.6 ete3
2 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/broccoli_env
3 conda install -c bioconda fasttree
4 conda install -c bioconda diamond=0.9.35
5 cd /data/SBCS-MartinDuranLab/03-Giacomo/src/
6 git clone https://github.com/rderelle/Broccoli
```

broccoli_v1.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/
3 #$ -o /data/scratch/btx654/
4 #$ -j y
5 #$ -pe smp 26
6 #$ -l h_vmem=20G
7 #$ -l h_rt=120:0:0
```

```
8  #$ -l highmem
9
10 module load anaconda3
11 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
    3/broccoli_env
12
13 mkdir broccoli
14 cd broccoli
15 cp -r /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/ ./
16
17 python /data/SBCS-MartinDuranLab/03-Giacomo/src/Broccoli/broccoli
    i.py -dir ./NR_proteomes -threads 26
```

- the fasta file in the folder “NR_proteomes” should all have the expansion .fasta (e.g. “Ofus.fasta”)

Broccoli -

```
1  module load anaconda3
2  conda create --prefix /data/SBCS-MartinDuranLab/03-Giacomo/src/a
    naconda3/broccoli_veryensitive_env python=3.6 ete3
3  conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
    3/broccoli_veryensitive_env
4  conda install -c bioconda fasttree
5  conda install -c bioconda diamond=2.0.6
6  cd /data/SBCS-MartinDuranLab/03-Giacomo/src/
7  cp -r Broccoli/ Broccoli_very_sensitive/
8  cp -r Broccoli/ Broccoli_ultra_sensitive/
9  #we need to modify the script Broccoli/scripts/broccoli_step2.py
    at line 238 changing --more-sensitive to --very-sensitive and --
    ultra-sensitive
```

broccoli_very_sensitive_v1.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
```

```
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 26
6  #$ -l h_vmem=20G
7  #$ -l h_rt=120:0:0
8  #$ -l highmem
9
10 module load anaconda3
11 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
   3/broccoli_veryensitive_env
12
13 mkdir -p broccoli_very_sensitive
14 cd broccoli_very_sensitive
15 cp -r /data/scratch/btx654/broccoli/NR_proteomes/ ./
16
17 python /data/SBCS-MartinDuranLab/03-Giacomo/src/Broccoli_very_se
   nsitive/broccoli.py -dir ./NR_proteomes -threads 26
```

broccoli_ultra_sensitive_v1.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 26
6  #$ -l h_vmem=20G
7  #$ -l h_rt=120:0:0
8  #$ -l highmem
9
10 module load anaconda3
11 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
   3/broccoli_veryensitive_env
```

```

12
13 mkdir -p broccoli_ultra_sensitive
14 cd broccoli_ultra_sensitive
15 cp -r /data/scratch/btx654/broccoli/NR_proteomes/ ./
16
17 python /data/SBCS-MartinDuranLab/03-Giacomo/src/Broccoli_ultra_s
    ensitive/broccoli.py -dir ./NR_proteomes -threads 26

```

OrthoFinder

/data/home/btx654/scripts/gene_family_evolution/orthofinder_Jan2021_v1.sh

```

1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 12
6  #$ -l h_vmem=5G
7  #$ -l h_rt=240:0:0
8
9
10 cd /data/scratch/btx654/gene_family_evolution/NR_proteomes/Ortho
    Finder/Results_Dec15/
11 #Load anaconda and activate MMseqs2 environment
12 module load anaconda3
13 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
    3/orthofinder_env
14
15 echo "(Nvec,(Hmia,(((Skow,Spur),(Blan,(Locu,Hsap))),((Smar,Tca
    s),(Smed,((Lgig,(Cgig,(Myes,Bpla))),(((Ofus,(Dgyr,(((Lluy,(Oalv,
    Rpac)),Ofra),(Ctel,(Hrob,Eand)))))),(Ngen,(Paus,Lana))))))));"
    > SpeciesTree_Jan2021.nwk
16 #Run orthofinder with mmseqs and inflation of 2

```



```
17 orthofinder -t 12 -S mmseqs -I 2 -s SpeciesTree_Jan2021.nwk -fg
/data/scratch/btx654/gene_family_evolution/NR_proteomes/OrthoFinder/Results_Dec15/
```

orthofinder_env

```
1 module load anaconda3
2 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
3/orthofinder_env
3 conda install -c bioconda orthofinder
• v 2.5.2-hdfd78af_1
```

/data/home/btx654/scripts/gene_family_evolution/Jun2021
/orthofinder_ultra_sensitive.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/
3 #$ -o /data/scratch/btx654/
4 #$ -j y
5 #$ -pe smp 12
6 #$ -l h_vmem=20G
7 #$ -l h_rt=240:0:0
8 #$ -l highmem
9
10 mkdir -p gene_family_evolution
11 cd gene_family_evolution
12 mkdir -p orthofinder_ultra_sensitive_Jun2021
13 cd orthofinder_ultra_sensitive_Jun2021
14 module load anaconda3
15 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
16 3/orthofinder_env
17
17 cp -r /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/ ./
18 echo "(Nvec,(Hmia,(((Skow,Spur),(Blan,(Locu,Hsap))),((Smar,Tca
```

```
s),(Smed,(((Lgig,Gaeg),(Cgig,(Myes,Bpla))),(((Ofus,(Dgyr,(((Llu
y,(Pech,(Oalv,Rpac))),Ofra),(Ctel,(Hrob,Eand)))))),(Ngen,(Paus,L
ana))))))));" > SpeciesTree.nwk
```

```
19 #Run orthofinder with mmseqs and inflation of 2
```

```
20 orthofinder -f NR_proteomes -t 12 -S diamond_ultra_sens -I 2 -s
SpeciesTree.nwk
```

/data/home/btx654/scripts/gene_family_evolution/Jun2021

/orthofinder_more_sensitive.sh

```
1 #!/bin/bash
```

```
2 #$ -wd /data/scratch/btx654/
```

```
3 #$ -o /data/scratch/btx654/
```

```
4 #$ -j y
```

```
5 #$ -pe smp 12
```

```
6 #$ -l h_vmem=20G
```

```
7 #$ -l h_rt=240:0:0
```

```
8 #$ -l highmem
```

```
9
```

```
10 mkdir -p gene_family_evolution
```

```
11 cd gene_family_evolution
```

```
12 mkdir -p orthofinder_more_sensitive_Jun2021
```

```
13 cd orthofinder_ultra_more_Jun2021
```

```
14 module load anaconda3
```

```
15 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
3/orthofinder_env
```

```
16
```

```
17 cp -r /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/ ./
```

```
18 echo "(Nvec,(Hmia,(((Skow,Spur),(Blan,(Locu,Hsap))),((Smar,Tca
s),(Smed,(((Lgig,Gaeg),(Cgig,(Myes,Bpla))),(((Ofus,(Dgyr,(((Llu
y,(Pech,(Oalv,Rpac))),Ofra),(Ctel,(Hrob,Eand)))))),(Ngen,(Paus,L
ana))))))));" > SpeciesTree.nwk
```

```
19 #Run orthofinder with mmseqs and inflation of 2
```

```
20 orthofinder -f NR_proteomes -t 12 -S diamond -I 2 -s SpeciesTree.nwk
```

Ferdi script

In Apple M1

```
1 conda install -c etetoolkit ete3
2 conda install numpy scipy statsmodels
3 conda install -c conda-forge tqdm
4 conda install -c anaconda ipykernel
5 python -m ipykernel install --user
6 conda install -c conda-forge notebook
7 conda install -c conda-forge ipywidgets
8 jupyter notebook # this will launch the browser app to actually
run the script
```

```
1 # first import
2 from collections import defaultdict
3 from collections import Counter
4 from ete3 import PhyloTree
5 import numpy as np
6 from tqdm import tqdm_notebook
7 import scipy.stats as stats
8 from statsmodels.stats.multitest import multipletests
9 import time
10 import csv
11 import copy
12
13 # broccoli
14 Fams,species={},[]
15 for line in open('table_0Gs_protein_names_modified_OFUS.txt'):
```

```
16     if line.startswith("#"): continue
17     fid=line.strip().split("\t")[0]
18     glt=[]
19     for i in line.strip().split("\t")[1:]:
20         glt.append(i)
21     gbs=defaultdict(list)
22     #print (fid)
23     for g in glt:
24         for j in g.split():
25             #if j =='gene-ND4L': continue
26             spc,gid=j.split('_',1)
27             #print (spc)
28             #print (gid)
29             gbs[spc].append(gid)
30     species.extend(gbs.keys())
31     Fams[fid]=gbs
32     #if i==5: break
33 species=list(set(species))
34 print(Fams['OG_1'])
35
36 # orthofinder
37 Fams={}
38 with open('Orthogroups.csv') as csvfile:
39     species=[]
40     for i,rc in enumerate(csv.reader(csvfile, delimiter='\t', qu
41 otechar='')):
42         #print (i)
43         #print (rc)
44         if i==0:
45             species=[sp for sp in rc[1:]]
```

```
45         continue
46     #for j,spg in enumerate(rc[1:]):
47         #for g in spg.split(','):
48             #print (g)
49             #print (g.split('|')[1])
50     #print (i)
51     #print (species)
52     #print(rc[1:])
53     #print(rc[19])
54     #print("hello")
55     genBySp=dict((species[j],[g.split('|')[1] for g in spg.s
56 plit(',')])) for j,spg in enumerate(rc[1:]) if not spg=='')
57     #print (genBySp)
58     #genBySp=
59     Fams[rc[0]]=genBySp
60     #if i==5: break
61
62 print(Fams['OG00000000'])
63
64 # PFAM domains
65 Fams={}
66 with open('Pfam_species_ferdi_OK.csv') as csvfile:
67     species=[]
68     for i,rc in enumerate(csv.reader(csvfile, delimiter='\t', qu
69 otechar='')):
70         #print (i)
71         #print (rc)
72         if i==0:
73             species=[sp for sp in rc[1:]]
74             continue
75     #for j,spg in enumerate(rc[1:]):
```

```
74         #for g in spg.split(','):
75             #print (g)
76             #print (g.split('|')[1])
77         #print (i)
78         #print (species)
79         #print(rc[1:])
80         #print(rc[19])
81         #print("hello")
82         genBySp=dict((species[j],[g.split('|')[1] for g in spg.s
plit(',')] for j,spg in enumerate(rc[1:]) if not spg=='')
83         #print (genBySp)
84         #genBySp=
85         Fams[rc[0]]=genBySp
86         #if i==5: break
87
88     print(Fams['PF00010.27'])
89
90     #Creates non-redundant species set
91     spSet=set()
92     for fid,gb in Fams.items():
93         for sp in gb.keys():
94             spSet.add(sp)
95     print(spSet)
96
97     #Load species tree
98     spT=PhyloTree(open('tree.tree').read(), format=1)
99
100     # Testing family expansions
101     #only testing species expansion
102     def fishExp(gbs,tbs):
```

```
103     cbs=dict((sp,len(gs)) for sp,gs in gbs.items())
104     med=np.median(list(cbs.values()))
105     #print(cbs,med)
106     spval={}
107     for sp,st in cbs.items():
108         nis=sum(cbs.values())-st
109         nit=sum(tbs.values())-tbs[sp]-sum(cbs.values())+st
110         odds,pval=stats.fisher_exact([[st, tbs[sp]-st], [nis, nit]],alternative='greater')
111         nmed=np.median([cbs[ssp] for ssp in cbs if not ssp==sp])
112         tmed=np.median([tbs[ssp] for ssp in tbs if not ssp==sp])
113         m_odds,m_pval=stats.fisher_exact([[st, tbs[sp]-st], [nmed, tmed]],alternative='greater')
114         #print(sp,pval,[st, tbs[sp]-st], [nis, nit],m_pval,[st, tbs[sp]-st], [nmed, tmed])
115         spval[sp]=(cbs[sp],m_pval)
116
117     return spval
118
119 #calculate number of genes per species
120 tbs=defaultdict(int)
121 for fid,gbs in Fams.items():
122     for sp in gbs:
123         tbs[sp]+=len(gbs[sp])
124
125 toosmall=0
126 expRes=[]
127 for fid,gbs in tqdm_notebook(Fams.items()):
128     cbs=dict((sp,len(gs)) for sp,gs in gbs.items())
129     if len(cbs.keys())<3 or sum(cbs.values())<5:
130         toosmall+=1
```

```
131         continue
132     try:
133         spval=fishExp(gbs,tbs)
134         expRes.append((fid,spval))
135     except:
136         print(fid,cbs)
137
138
139 print(len(Fams),len(expRes),toosmall)
140
141 spFam,spPval=defaultdict(list),defaultdict(list)
142 for fid,spval in expRes:
143     for sp in spval:
144         ct,pval=spval[sp]
145         spFam[sp].append(fid)
146         spPval[sp].append(pval)
147
148 famEnrich=defaultdict(dict)
149 for sp in spPval:
150     #signif_pvals = sidak(spPval[sp], alpha=0.05)
151     #before_adj=['True' if p < 0.05 else 'False' for p in spPval[sp] ]
152     #lsignif_pvals = lsu(np.array(spPval[sp]), q=0.05) #this is
benferroni-hochberg
153     adj=multiptests(pvals=np.array(spPval[sp]), alpha=0.05, method="fdr_bh")
154     print (sp,Counter(adj[0]))
155     for fam,pth,pval in zip(spFam[sp],adj[0],adj[1]):
156         famEnrich[fam][sp]=(pth,pval)
157
158 with open('orthogroups_ultrasensitive_enrich_Jun2021.txt','w') a
```



```
s out:
159     for fid, gbs in Fams.items():
160         cbs=dict((sp, len(gs)) for sp, gs in gbs.items())
161         pvl=famEnrich[fid]
162         enriched=set([])
163         outList=[fid,]
164         out.write('\t'.join(outList))
165 #Gains and Losses
166
167 def checkFam(fam, spTd):
168     fid, gbs=fam
169     spset=list(gbs.keys())
170     if len(spset)>1:
171         phtyp=spT.get_common_ancestor(list(gbs.keys()))
172     else:
173         phtyp=[l for l in spT.get_leaves() if l.name==spset
174 [0]][0]
175         #print phtyp.name
176         #phyloCt[phtyp.name]+=1
177         #gbs[Ai]
178         ndesc=len([l.name for l in phtyp.get_leaves()])
179         for leaf in spTd:
180             pv='1' if leaf.name in gbs else '0'
181             leaf.add_features(presence=pv)
182             lost=[node.name for node in phtyp.get_monophyletic(values=
183 ['0'], target_attr="presence")]
184             #print spTd.get_ascii(attributes=["name", "presence"], show_i
185 nternal=True)
186             #for node in spTd.get_monophyletic(values=['0'], target_attr
187 ="presence"):
188                 #    print node.name
```

```
185     # print node.get_ascii(attributes=["presence", "name"], s
how_internal=False)
186     return fid, len(gbs.keys()), ndesc, phtyp.name, lost
187
188 origins=defaultdict(int)
189 losses=defaultdict(int)
190 idGainLoss={}
191 for fam in Fams.items():
192     fid, nspec, ndesc, oritax, lost=checkFam(fam, spT)
193     idGainLoss[fid]=(nspec, ndesc, oritax, lost)
194     origins[oritax]+=1
195     for tax in lost:
196         losses[tax]+=1
197
198 spTi=copy.deepcopy(spT)
199 for node in spTi.traverse("postorder"):
200     #print node.name, origins[node.name], losses[node.name]
201     #node.add_features(origins=origins[node.name], losses=losses
[node.name])
202     node.name="{0}_{1}_{2}".format(node.name, origins[node.nam
e], losses[node.name])
203
204 print(spTi.write(format=7))
205 #print(spTi.get_ascii(attributes=["name"], show_internal=True))
206
207 #Output
208
209 with open('orthofinder_ultrasensitive_stats_Jun2021.tsv', 'w') as
out:
210     out.write("FID\tnb_genes\tnb_species\torigin\tlost_sp\tlost_
taxa\texpanded_tax\n")
```

```

211     for fid, gbs in Fams.items():
212         nspec, ndesc, oritax, lost = idGainLoss[fid]
213         ngn = sum([len(gl) for gl in gbs.values()])
214         exp = [s for s, p in famEnrich[fid].items() if p[0]]
215         #rpvals = zip(pvals.keys(), [round(-np.log10(l[1]), 2) for l
in pvals.values()])
216         outL = [fid, str(ngn), str(nspec), oritax, str(ndesc - nspe
c), ',', '.join(lost), ',', '.join(exp)]
217         out.write('\t'.join(outL) + '\n')

```

We need two input files for this:

- the .tsv output of OrthoFinder that we will need to transform into a csv (/orthofinder_Dec2020/Orthogroups/Orthogroups.tsv)
- and a newick format tree: tree.tree

```

(Nvec, (Hmia, (((Skow, Spur) Ambulacraria, (Blan, (Locu, Hsap) Vertebrat
a) Chordata) Deuterostomia, ((Smar, Tcas) Arthropoda, (Smed, (((Lgig, Ga
eg) Gastropoda, (Cgig, (Myes, Bpla) Bivalvia_cl1) Bivalvia) Mollusc
a, ((Ofus, (Dgyr, (((Lluy, (Pech, (Oalv, Rpac) Vestimentifera_cl2) Vesti
mentifera_cl1) Vestimentifera, Ofra) Siboglinidae, (Ctel, (Hrob, Eand)
Clitellata) Sedentaria_cl1) Sedentaria) Annelida_cl1) Annelida, (Nge
n, (Paus, Lana) Lophophorata) Kryptotrochozoa) Lophotrochozoa_cl1) Lop
hotrochozoa) Spiralia) Protostomia) Nephrozoa) Bilateria) Eumetazoa;

```

First thing first let's edit the Orthogroups.tsv file:

```

1 cp Orthogroups.tsv Orthogroups_original.tsv
2 sed -i 's/OFUS/Ofus|OFUS/g' Orthogroups.tsv
3 sed -i 's/Blan_/Blan|/g' Orthogroups.tsv
4 sed -i 's/Bpla_/Bpla|/g' Orthogroups.tsv

```

```
5 sed -i 's/Cgig_/Cgig|/g' Orthogroups.tsv
6 sed -i 's/Ctel_/Ctel|/g' Orthogroups.tsv
7 sed -i 's/Dgyr_/Dgyr|/g' Orthogroups.tsv
8 sed -i 's/Eand_/Eand|/g' Orthogroups.tsv
9 sed -i 's/Hmia_/Hmia|/g' Orthogroups.tsv
10 sed -i 's/Hrob_/Hrob|/g' Orthogroups.tsv
11 sed -i 's/Hsap_/Hsap|/g' Orthogroups.tsv
12 sed -i 's/Lana_/Lana|/g' Orthogroups.tsv
13 sed -i 's/Lgig_/Lgig|/g' Orthogroups.tsv
14 sed -i 's/Lluy_/Lluy|/g' Orthogroups.tsv
15 sed -i 's/Locu_/Locu|/g' Orthogroups.tsv
16 sed -i 's/Myes_/Myes|/g' Orthogroups.tsv
17 sed -i 's/Ngen_/Ngen|/g' Orthogroups.tsv
18 sed -i 's/Nvec_/Nvec|/g' Orthogroups.tsv
19 sed -i 's/Oalv_/Oalv|/g' Orthogroups.tsv
20 sed -i 's/Ofra_/Ofra|/g' Orthogroups.tsv
21 sed -i 's/Ofus_/Ofus|/g' Orthogroups.tsv
22 sed -i 's/Paus_/Paus|/g' Orthogroups.tsv
23 sed -i 's/Rpac_/Rpac|/g' Orthogroups.tsv
24 sed -i 's/Skow_/Skow|/g' Orthogroups.tsv
25 sed -i 's/Smar_/Smar|/g' Orthogroups.tsv
26 sed -i 's/Smed_/Smed|/g' Orthogroups.tsv
27 sed -i 's/Spur_/Spur|/g' Orthogroups.tsv
28 sed -i 's/Tcas_/Tcas|/g' Orthogroups.tsv
29 sed -i 's/Spur_/Spur|/g' Orthogroups.tsv
30 sed -i 's/Gaeg_/Gaeg|/g' Orthogroups.tsv
31 sed -i 's/Pech_/Pech|/g' Orthogroups.tsv
32
33 mv Orthogroups.tsv Orthogroups.csv
34 remove "orthogroups" from the first line with nano
```

```
35  
36 sed -i 's/gene-ND4L/Myes|gene-ND4L/g' Orthogroups.csv #this was  
creating a problem
```

Now we can launch the modified version of ferdi script and it will produce the results

```
1 cd /Users/giacomo/Jupyter_notebook/Ferdi_script/Jun2021/ultra_se  
nsitive  
2 cp /Users/giacomo/Jupyter_notebook/Ferdi_script/comp_genomics/Gi  
acomo.ipynb /Users/giacomo/Jupyter_notebook/Ferdi_script/Jun2021  
/ultra_sensitive  
3 conda activate Ferdi_env  
4 jupyter notebook  
5 cd /Users/giacomo/Jupyter_notebook/Ferdi_script/Jun2021/more_sen  
sitive  
6 cp /Users/giacomo/Jupyter_notebook/Ferdi_script/comp_genomics/Gi  
acomo.ipynb /Users/giacomo/Jupyter_notebook/Ferdi_script/Jun2021  
/more_sensitive
```

Expansion, Gains and Losses

Expansions

orthogroups_annotations_expanded_oasisia.sh

```
1 #!/bin/bash  
2 #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J  
ul2021/expansions/oasisia  
3 #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju  
l2021/expansions/oasisia  
4 #$ -j y  
5 #$ -pe smp 1  
6 #$ -l h_vmem=100G
```

```
7  # $ -l h_rt=72:00:0
8  # $ -l highmem
9
10 cut -f 1,7 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
    grep Oalv | cut -f 1 > gene_families_expanded_oasisia.txt #families
    expanded in oasisia
11 fgrep -f gene_families_expanded_oasisia.txt ../../Orthogroups.csv
    v > gene_families_expanded_oasisia.csv
12
13 cut -f 1,19 gene_families_expanded_oasisia.csv > orthogroups_gene_
    IDs_expanded_oasisia.txt
14 sed 's/Oalv|//g' orthogroups_gene_IDs_expanded_oasisia.txt > ortho
    groups_gene_IDs_expanded_oasisia_OK.txt
15
16 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t'
    '"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_number"
    > orthogroups_annotations_expanded_oasisia.csv
17 while read line; do
18     genes=$(cut -f 2 <<< "$line")
19     echo $genes
20     orthogroup_ID=$(cut -f 1 <<< "$line")
21     echo $orthogroup_ID
22     if [[ "$genes" == OALV* ]]
23     then
24         IFS=', '      # space is set as delimiter
25         read -ra ADDR <<< "$genes"    # str is read into an array a
            s tokens separated by IFS
26         for gene in "${ADDR[@]}; do
27             annotations=$(cut -f 13,14,15,18,20,24 ../../oasisia_anno
                tation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
28             echo $orthogroup_ID$'\t'"oasisia"$'\t'$annotations >> or
                thogroups_annotations_expanded_oasisia.csv
```

```

29     done
30     else
31         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
        ""$'\t'"" >> orthogroups_annotations_expanded_oasisia.csv
32     fi
33 done < orthogroups_gene_IDs_expanded_oasisia_OK.txt
34
35 echo "Orthogroup"$'\t'"Panther_annotation" > orthogroups_mostAbu
        ndantAnnotation_expanded_oasisia.csv
36 while read line; do
37     orthogroup_ID=$(cut -f 1 <<< "$line")
38     annotation=$(fgrep $orthogroup_ID orthogroups_annotatio
        ns_expanded_oasisia.csv | cut -f 7 | sort | uniq -c | sort -r | awk
        '{$1=""}; print $0}' | head -1)
39     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
        bundantAnnotation_expanded_oasisia.csv
40 done < gene_families_expanded_oasisia.csv

```

orthogroups_annotations_expanded_osedax.sh

```

1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
        ul2021/expansions/osedax
3  #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
        l2021/expansions/osedax
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=100G
7  #$ -l h_rt=72:00:0
8  #$ -l highmem
9
10 cut -f 1,7 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
    grep Ofra | cut -f 1 > gene_families_expanded_osedax.txt #famili

```

```

es expanded in oasisia
11 fgrep -f gene_families_expanded_osedax.txt ../../Orthogroups.csv
> gene_families_expanded_osedax.csv
12
13 cut -f 1,20 gene_families_expanded_osedax.csv > orthogroups_gene
_IDs_expanded_osedax.txt
14 sed 's/Ofra|//g' orthogroups_gene_IDs_expanded_osedax.txt > orth
ogroups_gene_IDs_expanded_osedax_OK.txt
15
16 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t'
\t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
mber" > orthogroups_annotations_expanded_osedax.csv
17 while read line; do
18     genes=$(cut -f 2 <<< "$line")
19     echo $genes
20     orthogroup_ID=$(cut -f 1 <<< "$line")
21     echo $orthogroup_ID
22     if [[ "$genes" == OFRA* ]]
23     then
24         IFS=', '      # space is set as delimiter
25         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
26         for gene in "${ADDR[@]}; do
27             annotations=$(cut -f 13,14,15,18,20,24 ../../osedax_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
28             echo $orthogroup_ID$'\t'"osedax"$'\t'"$annotations >> ort
hogroups_annotations_expanded_osedax.csv
29             done
30         else
31             echo $orthogroup_ID$'\t'"'"$'\t'"'"$'\t'"'"$'\t'"'"$'\t'"'"$'\t'
'"'"$'\t'"'" >> orthogroups_annotations_expanded_osedax.csv
32         fi

```



```
33 done < orthogroups_gene_IDs_expanded_osedax_OK.txt
34
35 echo "Orthogroup"${'\t'"Panther_annotation" > orthogroups_mostAbu
ndantAnnotation_expanded_osedax.csv
36 while read line; do
37     orthogroup_ID=$(cut -f 1 <<< "$line")
38     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_exp
anded_osedax.csv | cut -f 7 | sort | uniq -c | sort -r | awk
'{$1=""; print $0}' | head -1)
39     echo $orthogroup_ID${'\t'}$annotation >> orthogroups_mostA
bundantAnnotation_expanded_osedax.csv
40 done < gene_families_expanded_osedax.csv
```

orthogroups_annotations_expanded_riftia.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
ul2021/expansions/riftia
3 #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
l2021/expansions/riftia
4 #$ -j y
5 #$ -pe smp 1
6 #$ -l h_vmem=100G
7 #$ -l h_rt=72:00:0
8 #$ -l highmem
9
10 cut -f 1,7 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
grep Rpac | cut -f 1 > gene_families_expanded_riftia.txt #famili
es expanded in riftia
11 fgrep -f gene_families_expanded_riftia.txt ../../Orthogroups.csv
> gene_families_expanded_riftia.csv
12
```

```

13 cut -f 1,24 gene_families_expanded_riftia.csv > orthogroups_gene
   _IDs_expanded_riftia.txt
14 sed 's/Rpac|//g' orthogroups_gene_IDs_expanded_riftia.txt > orth
   ogroups_gene_IDs_expanded_riftia_OK.txt
15
16 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t'
   t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
   mber" > orthogroups_annotations_expanded_riftia.csv
17 while read line; do
18     genes=$(cut -f 2 <<< "$line")
19     echo $genes
20     orthogroup_ID=$(cut -f 1 <<< "$line")
21     echo $orthogroup_ID
22     if [[ "$genes" == RPAC* ]]
23     then
24         IFS=', '      # space is set as delimiter
25         read -ra ADDR <<< "$genes"    # str is read into an array a
   s tokens separated by IFS
26         for gene in "${ADDR[@]}"; do
27             annotations=$(cut -f 13,14,15,18,20,24 ../../riftia_annot
   ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
28             echo $orthogroup_ID$'\t'"riftia"$'\t'$annotations >> ort
   hogroups_annotations_expanded_riftia.csv
29         done
30     else
31         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
   ""$'\t'"" >> orthogroups_annotations_expanded_riftia.csv
32     fi
33 done < orthogroups_gene_IDs_expanded_riftia_OK.txt
34
35 echo "Orthogroup"$'\t'"Panther_annotation" > orthogroups_mostAbu
   ndantAnnotation_expanded_riftia.csv

```

```

36 while read line; do
37     orthogroup_ID=$(cut -f 1 <<< "$line")
38     annotation=$(fgrep $orthogroup_ID orthogroups_annotat
39     ions_expanded_riftia.csv | cut -f 7 | sort | uniq -c | sort -r | awk
40     '{ $1="" ; print $0 }' | head -1)
41     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
42     bundantAnnotation_expanded_riftia.csv
43 done < gene_families_expanded_riftia.csv

```

orthogroups_annotat ions_expanded_lamellibrachia.sh

```

1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
3  ul2021/expansions/lamellibrachia
4  #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
5  l2021/expansions/lamellibrachia
6  #$ -j y
7  #$ -pe smp 1
8  #$ -l h_vmem=100G
9  #$ -l h_rt=72:00:0
10 #$ -l highmem
11
12 cut -f 1,7 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
13 grep Lluy | cut -f 1 > gene_families_expanded_lamellibrachia.txt
14 #families expanded in riftia
15 fgrep -f gene_families_expanded_lamellibrachia.txt ../../Orthogr
16 oups.csv > gene_families_expanded_lamellibrachia.csv
17
18 cut -f 1,14 gene_families_expanded_lamellibrachia.csv > orthogro
19 ups_gene_IDs_expanded_lamellibrachia.txt
20 sed 's/Lluy|//g' orthogroups_gene_IDs_expanded_lamellibrachia.tx
21 t > orthogroups_gene_IDs_expanded_lamellibrachia_OK.txt
22
23 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t'

```

```

t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_number" > orthogroups_annotations_expanded_lamellibrachia.csv
17 while read line; do
18     genes=$(cut -f 2 <<< "$line")
19     echo $genes
20     orthogroup_ID=$(cut -f 1 <<< "$line")
21     echo $orthogroup_ID
22     if [[ "$genes" == FUN* ]]
23     then
24         IFS=', '      # space is set as delimiter
25         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
26         for gene in "${ADDR[@]}; do
27             annotations=$(cut -f 13,14,15,18,20,24 ../../lamellibrach
ia_annotation_Feb2021_TrinoPantherKO_OK.xls | fgrep $gene)
28             echo $orthogroup_ID$'\t'"lamellibrachia"$'\t'$annotation
s >> orthogroups_annotations_expanded_lamellibrachia.csv
29         done
30     else
31         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
""$'\t'"" >> orthogroups_annotations_expanded_lamellibrachia.cs
v
32     fi
33 done < orthogroups_gene_IDs_expanded_lamellibrachia_OK.txt
34
35 echo "Orthogroup"$'\t'"Panther_annotation" > orthogroups_mostAbu
ndantAnnotation_expanded_lamellibrachia.csv
36 while read line; do
37     orthogroup_ID=$(cut -f 1 <<< "$line")
38     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_exp
anded_lamellibrachia.csv | cut -f 7 | sort | uniq -c | sort -r
| awk '{ $1="" ; print $0 }' | head -1)

```

```
39     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
    abundantAnnotation_expanded_lamellibrachia.csv
40 done < gene_families_expanded_lamellibrachia.csv
```

orthogroups_annotations_expanded_paraescarpia.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
    ul2021/expansions/paraescarpia
3  #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
    l2021/expansions/paraescarpia
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=100G
7  #$ -l h_rt=72:00:0
8  #$ -l highmem
9
10 cut -f 1,7 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
    grep Pech | cut -f 1 > gene_families_expanded_paraescarpia.txt #
    families expanded in riftia
11 fgrep -f gene_families_expanded_paraescarpia.txt ../../Orthogrou
    ps.csv > gene_families_expanded_paraescarpia.csv
12
13 cut -f 1,23 gene_families_expanded_paraescarpia.csv > orthogroup
    s_gene_IDs_expanded_paraescarpia.txt
14 sed 's/Pech|//g' orthogroups_gene_IDs_expanded_paraescarpia.txt
    > orthogroups_gene_IDs_expanded_paraescarpia_OK.txt
15
16 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t'
    "'GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
    mber" > orthogroups_annotations_expanded_paraescarpia.csv
17 while read line; do
18     genes=$(cut -f 2 <<< "$line")
```

```

19  echo $genes
20  orthogroup_ID=$(cut -f 1 <<< "$line")
21  echo $orthogroup_ID
22  if [[ "$genes" == nbis* ]]
23  then
24      IFS=', '      # space is set as delimiter
25      read -ra ADDR <<< "$genes"      # str is read into an array a
s tokens separated by IFS
26      for gene in "${ADDR[@]}"; do
27          annotations=$(cut -f 13,14,15,18,20,24 ../../paraescarpia
_annotation_Jun2021_TrinoPantherK0.xls | fgrep $gene)
28          echo $orthogroup_ID$'\t'"paraescarpia"$'\t'$annotations
>> orthogroups_annotations_expanded_paraescarpia.csv
29      done
30  else
31      echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
""$'\t'"" >> orthogroups_annotations_expanded_paraescarpia.csv
32  fi
33 done < orthogroups_gene_IDs_expanded_paraescarpia_OK.txt
34
35 echo "Orthogroup"$'\t'"Panther_annotation" > orthogroups_mostAbu
ndantAnnotation_expanded_paraescarpia.csv
36 while read line; do
37     orthogroup_ID=$(cut -f 1 <<< "$line")
38     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_exp
anded_paraescarpia.csv | cut -f 7 | sort | uniq -c | sort -r |
awk '{ $1="" ; print $0 }' | head -1)
39     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
bundantAnnotation_expanded_paraescarpia.csv
40 done < gene_families_expanded_paraescarpia.csv

```

Gains

orthogroups_annotations_originated_oasisia.sh

```

1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
   ul2021/gains/oasisia
3  #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
   l2021/gains/oasisia
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=100G
7  #$ -l h_rt=72:00:0
8  #$ -l highmem
9
10 cut -f 1,4 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
   grep -w Oalv | cut -f 1 > gene_families_originated_oasisia.txt #
   families originated in oasisia
11 fgrep -f gene_families_originated_oasisia.txt ../../Orthogroups.
   csv > gene_families_originated_oasisia.csv
12
13 cut -f 1,19 gene_families_originated_oasisia.csv > orthogroups_g
   ene_IDs_originated_oasisia.txt
14 sed 's/Oalv|//g' orthogroups_gene_IDs_originated_oasisia.txt > o
   rthogroups_gene_IDs_originated_oasisia_OK.txt
15
16 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t'
   t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
   mber" > orthogroups_annotations_originated_oasisia.csv
17 while read line; do
18     genes=$(cut -f 2 <<< "$line")
19     echo $genes
20     orthogroup_ID=$(cut -f 1 <<< "$line")
21     echo $orthogroup_ID
22     if [[ "$genes" == OALV* ]]

```

```

23     then
24         IFS=', '      # space is set as delimiter
25         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
26         for gene in "${ADDR[@]}"; do
27             annotations=$(cut -f 13,14,15,18,20,24 ../../oasisia_anno
tation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
28             echo $orthogroup_ID$'\t'"oasisia"$'\t'$annotations >> or
thogroups_annotations_originated_oasisia.csv
29         done
30     else
31         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_originated_oasisia.csv
32     fi
33 done < orthogroups_gene_IDs_originated_oasisia_OK.txt
34
35 echo "Orthogroup"$'\t'"Panther_annotation" > orthogroups_mostAbu
ndantAnnotation_originated_oasisia.csv
36 while read line; do
37     orthogroup_ID=$(cut -f 1 <<< "$line")
38     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_ori
ginated_oasisia.csv | cut -f 7 | sort | uniq -c | sort -r | awk
'{$1=""; print $0}' | head -1)
39     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
bundantAnnotation_originated_oasisia.csv
40 done < gene_families_originated_oasisia.csv

```

orthogroups_annotations_originated_osedax.sh

```

1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
ul2021/gains/osedax

```



```
3  #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
   l2021/gains/osedax
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=100G
7  #$ -l h_rt=72:00:0
8  #$ -l highmem
9
10 cut -f 1,4 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
   grep -w Ofra | cut -f 1 > gene_families_originated_osedax.txt #f
   amilies originated in oasisia
11 fgrep -f gene_families_originated_osedax.txt ../../Orthogroups.c
   sv > gene_families_originated_osedax.csv
12
13 cut -f 1,20 gene_families_originated_osedax.csv > orthogroups_ge
   ne_IDs_originated_osedax.txt
14 sed 's/Ofra|//g' orthogroups_gene_IDs_originated_osedax.txt > or
   thogroups_gene_IDs_originated_osedax_OK.txt
15
16 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t
   t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
   mber" > orthogroups_annotations_originated_osedax.csv
17 while read line; do
18     genes=$(cut -f 2 <<< "$line")
19     echo $genes
20     orthogroup_ID=$(cut -f 1 <<< "$line")
21     echo $orthogroup_ID
22     if [[ "$genes" == OFRA* ]]
23     then
24         IFS=', '      # space is set as delimiter
25         read -ra ADDR <<< "$genes"    # str is read into an array a
   s tokens separated by IFS
```

```

26     for gene in "${ADDR[@]}; do
27         annotations=$(cut -f 13,14,15,18,20,24 ../../osedax_annotation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
28         echo $orthogroup_ID$'\t'"osedax"$'\t'$annotations >> orthogroups_annotations_originated_osedax.csv
29     done
30 else
31     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'"" >> orthogroups_annotations_originated_osedax.csv
32 fi
33 done < orthogroups_gene_IDs_originated_osedax_OK.txt
34
35 echo "Orthogroup"$'\t'"Panther_annotation" > orthogroups_mostAbundantAnnotation_originated_osedax.csv
36 while read line; do
37     orthogroup_ID=$(cut -f 1 <<< "$line")
38     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_originated_osedax.csv | cut -f 7 | sort | uniq -c | sort -r | awk '{ $1=""; print $0 }' | head -1)
39     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostAbundantAnnotation_originated_osedax.csv
40 done < gene_families_originated_osedax.csv

```

orthogroups_annotations_originated_riftia.sh

```

1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/Jul2021/gains/riftia
3 #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Jul2021/gains/riftia
4 #$ -j y
5 #$ -pe smp 1
6 #$ -l h_vmem=100G

```

```
7  # $ -l h_rt=72:00:0
8  # $ -l highmem
9
10 cut -f 1,4 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
    grep -w Rpac | cut -f 1 > gene_families_originated_riftia.txt #f
    amilies originated in riftia
11 fgrep -f gene_families_originated_riftia.txt ../../Orthogroups.c
    sv > gene_families_originated_riftia.csv
12
13 cut -f 1,24 gene_families_originated_riftia.csv > orthogroups_ge
    ne_IDs_originated_riftia.txt
14 sed 's/Rpac|//g' orthogroups_gene_IDs_originated_riftia.txt > or
    thogroups_gene_IDs_originated_riftia_OK.txt
15
16 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\
    t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
    mber" > orthogroups_annotations_originated_riftia.csv
17 while read line; do
18     genes=$(cut -f 2 <<< "$line")
19     echo $genes
20     orthogroup_ID=$(cut -f 1 <<< "$line")
21     echo $orthogroup_ID
22     if [[ "$genes" == RPAC* ]]
23     then
24         IFS=', '      # space is set as delimiter
25         read -ra ADDR <<< "$genes"    # str is read into an array a
            s tokens separated by IFS
26         for gene in "${ADDR[@]}"; do
27             annotations=$(cut -f 13,14,15,18,20,24 ../../riftia_annot
                ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
28             echo $orthogroup_ID$'\t'"riftia"$'\t'"$annotations >> ort
                hogroups_annotations_originated_riftia.csv
```

```

29     done
30     else
31         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
        ""$'\t'"" >> orthogroups_annotations_originated_riftia.csv
32     fi
33 done < orthogroups_gene_IDs_originated_riftia_OK.txt
34
35 echo "Orthogroup"$'\t'"Panther_annotation" > orthogroups_mostAbu
        ndantAnnotation_originated_riftia.csv
36 while read line; do
37     orthogroup_ID=$(cut -f 1 <<< "$line")
38     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_ori
        genated_riftia.csv | cut -f 7 | sort | uniq -c | sort -r | awk
        '{$1=""; print $0}' | head -1)
39     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
        bundantAnnotation_originated_riftia.csv
40 done < gene_families_originated_riftia.csv

```

orthogroups_annotations_originated_lamellibrachia.sh

```

1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
        ul2021/gains/lamellibrachia
3  #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
        l2021/gains/lamellibrachia
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=100G
7  #$ -l h_rt=72:00:0
8  #$ -l highmem
9
10 cut -f 1,4 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
        grep -w Lluy | cut -f 1 > gene_families_originated_lamellibrachi

```

```

a.txt #families originated in riftia
11 fgrep -f gene_families_originated_lamellibrachia.txt ../../Ortho
    groups.csv > gene_families_originated_lamellibrachia.csv
12
13 cut -f 1,14 gene_families_originated_lamellibrachia.csv > orthog
    rroups_gene_IDs_originated_lamellibrachia.txt
14 sed 's/Lluy|//g' orthogroups_gene_IDs_originated_lamellibrachia.
    txt > orthogroups_gene_IDs_originated_lamellibrachia_OK.txt
15
16 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t'
    t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
    mber" > orthogroups_annotations_originated_lamellibrachia.csv
17 while read line; do
18     genes=$(cut -f 2 <<< "$line")
19     echo $genes
20     orthogroup_ID=$(cut -f 1 <<< "$line")
21     echo $orthogroup_ID
22     if [[ "$genes" == FUN* ]]
23     then
24         IFS=', '      # space is set as delimiter
25         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
26         for gene in "${ADDR[@]}; do
27             annotations=$(cut -f 13,14,15,18,20,24 ../../lamellibrach
ia_annotation_Feb2021_TrinoPantherKO_OK.xls | fgrep $gene)
28             echo $orthogroup_ID$'\t'"lamellibrachia"$'\t'"$annotation
s >> orthogroups_annotations_originated_lamellibrachia.csv
29         done
30     else
31         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
        ""$'\t'"" >> orthogroups_annotations_originated_lamellibrachia.
        csv

```

```
32     fi
33 done < orthogroups_gene_IDs_originated_lamellibrachia_OK.txt
34
35 echo "Orthogroup"${'\t'"Panther_annotation" > orthogroups_mostAbu
ndantAnnotation_originated_lamellibrachia.csv
36 while read line; do
37     orthogroup_ID=$(cut -f 1 <<< "$line")
38     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_ori
ginated_lamellibrachia.csv | cut -f 7 | sort | uniq -c | sort -
r | awk '{ $1="" ; print $0 }' | head -1)
39     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
bundantAnnotation_originated_lamellibrachia.csv
40 done < gene_families_originated_lamellibrachia.csv
```

orthogroups_annotations_originated_paraescarpia.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
ul2021/gains/paraescarpia
3 #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
l2021/gains/paraescarpia
4 #$ -j y
5 #$ -pe smp 1
6 #$ -l h_vmem=100G
7 #$ -l h_rt=72:00:0
8 #$ -l highmem
9
10 cut -f 1,4 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
grep -w Pech | cut -f 1 > gene_families_originated_paraescarpia.
txt #families originated in riftia
11 fgrep -f gene_families_originated_paraescarpia.txt ../../Orthogr
oups.csv > gene_families_originated_paraescarpia.csv
12
```

```

13 cut -f 1,23 gene_families_originated_paraescarpia.csv > orthogro
ups_gene_IDs_originated_paraescarpia.txt
14 sed 's/Pech|//g' orthogroups_gene_IDs_originated_paraescarpia.tx
t > orthogroups_gene_IDs_originated_paraescarpia_OK.txt
15
16 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\
\t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
mber" > orthogroups_annotations_originated_paraescarpia.csv
17 while read line; do
18     genes=$(cut -f 2 <<< "$line")
19     echo $genes
20     orthogroup_ID=$(cut -f 1 <<< "$line")
21     echo $orthogroup_ID
22     if [[ "$genes" == nbis* ]]
23     then
24         IFS=', '      # space is set as delimiter
25         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
26         for gene in "${ADDR[@]}"; do
27             annotations=$(cut -f 13,14,15,18,20,24 ../../paraescarpia
_annotation_Jun2021_TrinoPantherK0.xls | fgrep $gene)
28             echo $orthogroup_ID$'\t'"paraescarpia"$'\t'"$annotations
>> orthogroups_annotations_originated_paraescarpia.csv
29         done
30     else
31         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_originated_paraescarpia.cs
v
32     fi
33 done < orthogroups_gene_IDs_originated_paraescarpia_OK.txt
34
35 echo "Orthogroup"$'\t'"Panther_annotation" > orthogroups_mostAbu

```

```
ndantAnnotation_originated_paraescarpia.csv
36 while read line; do
37     orthogroup_ID=$(cut -f 1 <<< "$line")
38     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_ori
        ginated_paraescarpia.csv | cut -f 7 | sort | uniq -c | sort -r
        | awk '{ $1="" ; print $0 }' | head -1)
39     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
        bundantAnnotation_originated_paraescarpia.csv
40 done < gene_families_originated_paraescarpia.csv
```

orthogroups_annotations_originated_siboglinidae.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
    ul2021/gains/siboglinidae
3 #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
    l2021/gains/siboglinidae
4 #$ -j y
5 #$ -pe smp 1
6 #$ -l h_vmem=100G
7 #$ -l h_rt=140:00:0
8 #$ -l highmem
9
10 cut -f 1,4 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
    grep -w Siboglinidae | cut -f 1 > gene_families_originated_sibog
    linidae.txt #families originated in oasisia
11 fgrep -f gene_families_originated_siboglinidae.txt ../../Orthogr
    oups.csv > gene_families_originated_siboglinidae.csv
12
13 cut -f 1,19 gene_families_originated_siboglinidae.csv > orthogro
    ups_gene_IDs_originated_siboglinidae_oasisia.txt
14 sed 's/Oalv|//g' orthogroups_gene_IDs_originated_siboglinidae_oa
    sisia.txt > orthogroups_gene_IDs_originated_siboglinidae_oasisia
    _OK.txt
```



```
15 cut -f 1,20 gene_families_originated_siboglinidae.csv > orthogroups_gene_IDs_originated_siboglinidae_osedax.txt
16 sed 's/Ofra|//g' orthogroups_gene_IDs_originated_siboglinidae_osedax.txt > orthogroups_gene_IDs_originated_siboglinidae_osedax_OK.txt
17 cut -f 1,24 gene_families_originated_siboglinidae.csv > orthogroups_gene_IDs_originated_siboglinidae_riftia.txt
18 sed 's/Rpac|//g' orthogroups_gene_IDs_originated_siboglinidae_riftia.txt > orthogroups_gene_IDs_originated_siboglinidae_riftia_OK.txt
19 cut -f 1,14 gene_families_originated_siboglinidae.csv > orthogroups_gene_IDs_originated_siboglinidae_lamellibrachia.txt
20 sed 's/Lluy|//g' orthogroups_gene_IDs_originated_siboglinidae_lamellibrachia.txt > orthogroups_gene_IDs_originated_siboglinidae_lamellibrachia_OK.txt
21 cut -f 1,23 gene_families_originated_siboglinidae.csv > orthogroups_gene_IDs_originated_siboglinidae_paraescarpia.txt
22 sed 's/Pech|//g' orthogroups_gene_IDs_originated_siboglinidae_paraescarpia.txt > orthogroups_gene_IDs_originated_siboglinidae_paraescarpia_OK.txt
23
24
25 while read line; do
26     genes=$(cut -f 2 <<< "$line")
27     echo $genes
28     orthogroup_ID=$(cut -f 1 <<< "$line")
29     echo $orthogroup_ID
30     if [[ "$genes" == OALV* ]]
31     then
32         IFS=', '      # space is set as delimiter
33         read -ra ADDR <<< "$genes"    # str is read into an array as tokens separated by IFS
34         for gene in "${ADDR[@]}; do
```

```

35     annotations=$(cut -f 13,14,15,18,20,24 ../../oasisia_anno
tation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
36     echo $orthogroup_ID$'\t'"oasisia"$'\t'$annotations >> or
thogroups_annotations_originated_siboglinidae_oasisia.csv
37     done
38     else
39     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_originated_siboglinidae_oa
sisia.csv
40     fi
41 done < orthogroups_gene_IDs_originated_siboglinidae_oasisia_OK.t
xt
42
43 while read line; do
44     genes=$(cut -f 2 <<< "$line")
45     echo $genes
46     orthogroup_ID=$(cut -f 1 <<< "$line")
47     echo $orthogroup_ID
48     if [[ "$genes" == OFRA* ]]
49     then
50         IFS=', '      # space is set as delimiter
51         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
52         for gene in "${ADDR[@]}"; do
53             annotations=$(cut -f 13,14,15,18,20,24 ../../osedax_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
54             echo $orthogroup_ID$'\t'"osedax"$'\t'$annotations >> ort
hogroups_annotations_originated_siboglinidae_osedax.csv
55             done
56             else
57             echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_originated_siboglinidae_os

```

```
edax.csv
```

```
58     fi
59 done < orthogroups_gene_IDs_originated_siboglinidae_osedax_OK.txt
60
61 while read line; do
62     genes=$(cut -f 2 <<< "$line")
63     echo $genes
64     orthogroup_ID=$(cut -f 1 <<< "$line")
65     echo $orthogroup_ID
66     if [[ "$genes" == RPAC* ]]
67     then
68         IFS=', '      # space is set as delimiter
69         read -ra ADDR <<< "$genes"  # str is read into an array as
s tokens separated by IFS
70         for gene in "${ADDR[@]}; do
71             annotations=$(cut -f 13,14,15,18,20,24 ../../riftia_annotation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
72             echo $orthogroup_ID$'\t'"riftia"$'\t'$annotations >> orthogroups_annotations_originated_siboglinidae_riftia.csv
73         done
74     else
75         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'"" >> orthogroups_annotations_originated_siboglinidae_riftia.csv
76     fi
77 done < orthogroups_gene_IDs_originated_siboglinidae_riftia_OK.txt
78
79 while read line; do
80     genes=$(cut -f 2 <<< "$line")
81     echo $genes
```

```
82 orthogroup_ID=$(cut -f 1 <<< "$line")
83 echo $orthogroup_ID
84 if [[ "$genes" == FUN* ]]
85 then
86     IFS=', '      # space is set as delimiter
87     read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
88     for gene in "${ADDR[@]}; do
89         annotations=$(cut -f 13,14,15,18,20,24 ../../lamellibrach
ia_annotation_Feb2021_TrinoPantherKO_OK.xls | fgrep $gene)
90         echo $orthogroup_ID$'\t'"lamellibrachia"$'\t'$annotation
s >> orthogroups_annotations_originated_siboglinidae_lamellibrac
hia.csv
91     done
92 else
93     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_originated_siboglinidae_la
mellibrachia.csv
94 fi
95 done < orthogroups_gene_IDs_originated_siboglinidae_lamellibrach
ia_OK.txt
96
97 while read line; do
98     genes=$(cut -f 2 <<< "$line")
99     echo $genes
100     orthogroup_ID=$(cut -f 1 <<< "$line")
101     echo $orthogroup_ID
102     if [[ "$genes" == nbis* ]]
103     then
104         IFS=', '      # space is set as delimiter
105         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
```

```

106     for gene in "${ADDR[@]"}; do
107         annotations=$(cut -f 13,14,15,18,20,24 ../../paraescarpia
_annotation_Jun2021_TrinoPantherK0.xls | fgrep $gene)
108         echo $orthogroup_ID$'\t'"paraescarpia"$'\t'$annotations
>> orthogroups_annotations_originated_siboglinidae_paraescarpia.
csv
109     done
110     else
111         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
'""$'\t'"" >> orthogroups_annotations_originated_siboglinidae_pa
raescarpia.csv
112     fi
113 done < orthogroups_gene_IDs_originated_siboglinidae_paraescarpia
_OK.txt
114
115 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t'
t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
mber" > orthogroups_annotations_originated_siboglinidae_ofra_0al
v_Rpac_Lluy_Pech.csv
116 cat orthogroups_annotations_originated_siboglinidae_oasisia.csv
orthogroups_annotations_originated_siboglinidae_osedax.csv ortho
groups_annotations_originated_siboglinidae_riftia.csv orthogroup
s_annotations_originated_siboglinidae_lamellibrachia.csv orthogr
oups_annotations_originated_siboglinidae_paraescarpia.csv >> ort
hogroups_annotations_originated_siboglinidae_ofra_0alv_Rpac_Lluy
_Pech.csv
117 sort orthogroups_annotations_originated_siboglinidae_ofra_0alv_R
pac_Lluy_Pech.csv > orthogroups_annotations_originated_siboglini
dae_ofra_0alv_Rpac_Lluy_Pech_OK.csv
118
119 echo "Orthogroup"$'\t'"Panther_annotation" > orthogroups_mostAbu
ndantAnnotation_originated_siboglinidae.csv
120 while read line; do
121     orthogroup_ID=$(cut -f 1 <<< "$line")

```

```

122     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_ori
    genated_siboglinidae_0fra_0alv_Rpac_Lluy_Pech_OK.csv | cut -f 7
    | sed '/^$/d' | sort | uniq -c | sort -r | awk '{s1=""; print
    $0}' | head -1)
123         echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
    bundantAnnotation_originated_siboglinidae.csv
124 done < gene_families_originated_siboglinidae.csv

```

orthogroups_annotations_originated_vestimentifera.sh

```

1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
    ul2021/gains/vestimentifera
3  #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
    l2021/gains/vestimentifera
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=100G
7  #$ -l h_rt=140:00:0
8  #$ -l highmem
9
10 cut -f 1,4 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
    grep -w Vestimentifera | cut -f 1 > gene_families_originated_ves
    timentifera.txt #families originated in oasisia
11 fgrep -f gene_families_originated_vestimentifera.txt ../../Ortho
    groups.csv > gene_families_originated_vestimentifera.csv
12
13 cut -f 1,19 gene_families_originated_vestimentifera.csv > orthog
    rroups_gene_IDs_originated_vestimentifera_oasisia.txt
14 sed 's/0alv|//g' orthogroups_gene_IDs_originated_vestimentifera_
    oasisia.txt > orthogroups_gene_IDs_originated_vestimentifera_oas
    isia_OK.txt
15 cut -f 1,24 gene_families_originated_vestimentifera.csv > orthog
    rroups_gene_IDs_originated_vestimentifera_riftia.txt

```

```
16 sed 's/Rpac|//g' orthogroups_gene_IDs_originated_vestimentifera_
    riftia.txt > orthogroups_gene_IDs_originated_vestimentifera_rift
    ia_OK.txt
17 cut -f 1,14 gene_families_originated_vestimentifera.csv > orthog
    rroups_gene_IDs_originated_vestimentifera_lamellibrachia.txt
18 sed 's/Lluy|//g' orthogroups_gene_IDs_originated_vestimentifera_
    lamellibrachia.txt > orthogroups_gene_IDs_originated_vestimentif
    era_lamellibrachia_OK.txt
19 cut -f 1,23 gene_families_originated_vestimentifera.csv > orthog
    rroups_gene_IDs_originated_vestimentifera_paraescarpia.txt
20 sed 's/Pech|//g' orthogroups_gene_IDs_originated_vestimentifera_
    paraescarpia.txt > orthogroups_gene_IDs_originated_vestimentifer
    a_paraescarpia_OK.txt
21
22
23 while read line; do
24     genes=$(cut -f 2 <<< "$line")
25     echo $genes
26     orthogroup_ID=$(cut -f 1 <<< "$line")
27     echo $orthogroup_ID
28     if [[ "$genes" == OALV* ]]
29     then
30         IFS=', '      # space is set as delimiter
31         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
32         for gene in "${ADDR[@]}; do
33             annotations=$(cut -f 13,14,15,18,20,24 ../../oasisia_anno
tation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
34             echo $orthogroup_ID$'\t'"oasisia"$'\t'$annotations >> or
thogroups_annotations_originated_vestimentifera_oasisia.csv
35         done
36     else
```

```
37     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
    ""$'\t'"" >> orthogroups_annotations_originated_vestimentifera_
oasisia.csv
38     fi
39 done < orthogroups_gene_IDs_originated_vestimentifera_oasisia_0
K.txt
40
41 while read line; do
42     genes=$(cut -f 2 <<< "$line")
43     echo $genes
44     orthogroup_ID=$(cut -f 1 <<< "$line")
45     echo $orthogroup_ID
46     if [[ "$genes" == RPAC* ]]
47     then
48         IFS=', '      # space is set as delimiter
49         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
50         for gene in "${ADDR[@]}"; do
51             annotations=$(cut -f 13,14,15,18,20,24 ../../riftia_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
52             echo $orthogroup_ID$'\t'""riftia"$'\t'$annotations >> ort
hogroups_annotations_originated_vestimentifera_riftia.csv
53         done
54     else
55         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
    ""$'\t'"" >> orthogroups_annotations_originated_vestimentifera_
riftia.csv
56     fi
57 done < orthogroups_gene_IDs_originated_vestimentifera_riftia_0K.
txt
58
59 while read line; do
```



```
60 genes=$(cut -f 2 <<< "$line")
61 echo $genes
62 orthogroup_ID=$(cut -f 1 <<< "$line")
63 echo $orthogroup_ID
64 if [[ "$genes" == FUN* ]]
65 then
66     IFS=', '      # space is set as delimiter
67     read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
68     for gene in "${ADDR[@]}; do
69         annotations=$(cut -f 13,14,15,18,20,24 ../../lamellibrach
ia_annotation_Feb2021_TrinoPantherKO_OK.xls | fgrep $gene)
70         echo $orthogroup_ID$'\t'"lamellibrachia"$'\t'$annotation
s >> orthogroups_annotations_originated_vestimentifera_lamellibr
achia.csv
71     done
72 else
73     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_originated_vestimentifera_
lamellibrachia.csv
74 fi
75 done < orthogroups_gene_IDs_originated_vestimentifera_lamellibra
chia_OK.txt
76
77 while read line; do
78     genes=$(cut -f 2 <<< "$line")
79     echo $genes
80     orthogroup_ID=$(cut -f 1 <<< "$line")
81     echo $orthogroup_ID
82     if [[ "$genes" == nbis* ]]
83     then
84         IFS=', '      # space is set as delimiter
```

```

85     read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
86     for gene in "${ADDR[@]}; do
87         annotations=$(cut -f 13,14,15,18,20,24 ../../paraescarpia
_annotation_Jun2021_TrinoPantherK0.xls | fgrep $gene)
88         echo $orthogroup_ID$'\t'"paraescarpia"$'\t'$annotations
>> orthogroups_annotations_originated_vestimentifera_paraescarpia.csv
89     done
90     else
91         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
""$'\t'"" >> orthogroups_annotations_originated_vestimentifera_
paraescarpia.csv
92     fi
93 done < orthogroups_gene_IDs_originated_vestimentifera_paraescarpia_OK.txt
94
95 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t'
t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
mber" > orthogroups_annotations_originated_vestimentifera_Oalv_R
pac_Lluy_Pech.csv
96 cat orthogroups_annotations_originated_vestimentifera_oasisia.cs
v orthogroups_annotations_originated_vestimentifera_riftia.csv o
rthogroups_annotations_originated_vestimentifera_lamellibrachia.
csv orthogroups_annotations_originated_vestimentifera_paraescarpia.csv >> orthogroups_annotations_originated_vestimentifera_Oalv_Rpac_Lluy_Pech.csv
97 sort orthogroups_annotations_originated_vestimentifera_Oalv_Rpac_Lluy_Pech.csv > orthogroups_annotations_originated_vestimentifera_Oalv_Rpac_Lluy_Pech_OK.csv
98
99 echo "Orthogroup"$'\t'"Panther_annotation" > orthogroups_mostAbundantAnnotation_originated_vestimentifera.csv
100 while read line; do

```

```
101     orthogroup_ID=$(cut -f 1 <<< "$line")
102     annotation=$(fgrep $orthogroup_ID orthogroups_annotat
    originated_vestimentifera_0fra_0alv_Rpac_Lluy_Pech_OK.csv | cut -f
    7 | sed '/^$/d' | sort | uniq -c | sort -r | awk '{ $1="" ; print
    $0 }' | head -1)
103     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
    bundantAnnotation_originated_vestimentifera.csv
104 done < gene_families_originated_vestimentifera.csv
```

orthogroups_annotatations_originated_vestimentifera_cl1.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
    ul2021/gains/vestimentifera_cl1
3  #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
    l2021/gains/vestimentifera_cl1
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=100G
7  #$ -l h_rt=140:00:0
8  #$ -l highmem
9
10 cut -f 1,4 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
    grep -w Vestimentifera_cl1 | cut -f 1 > gene_families_originated
    _vestimentifera_cl1.txt #families originated in oasisia
11 fgrep -f gene_families_originated_vestimentifera_cl1.txt ../../O
    rthogroups.csv > gene_families_originated_vestimentifera_cl1.csv
12
13 cut -f 1,19 gene_families_originated_vestimentifera_cl1.csv > or
    thogroups_gene_IDs_originated_vestimentifera_cl1_oasisia.txt
14 sed 's/0alv|//g' orthogroups_gene_IDs_originated_vestimentifera
    _cl1_oasisia.txt > orthogroups_gene_IDs_originated_vestimentifera
    _cl1_oasisia_OK.txt
15 cut -f 1,24 gene_families_originated_vestimentifera_cl1.csv > or
```

```

thogroups_gene_IDs_originated_vestimentifera_cl1_riftia.txt
16 sed 's/Rpac|//g' orthogroups_gene_IDs_originated_vestimentifera_
   cl1_riftia.txt > orthogroups_gene_IDs_originated_vestimentifera_
   cl1_riftia_OK.txt
17 cut -f 1,23 gene_families_originated_vestimentifera_cl1.csv > or
   thogroups_gene_IDs_originated_vestimentifera_cl1_paraescarpia.tx
   t
18 sed 's/Pech|//g' orthogroups_gene_IDs_originated_vestimentifera_
   cl1_paraescarpia.txt > orthogroups_gene_IDs_originated_vestiment
   ifera_cl1_paraescarpia_OK.txt
19
20
21 while read line; do
22     genes=$(cut -f 2 <<< "$line")
23     echo $genes
24     orthogroup_ID=$(cut -f 1 <<< "$line")
25     echo $orthogroup_ID
26     if [[ "$genes" == OALV* ]]
27     then
28         IFS=', '      # space is set as delimiter
29         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
30         for gene in "${ADDR[@]}; do
31             annotations=$(cut -f 13,14,15,18,20,24 ../../oasisia_anno
tation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
32             echo $orthogroup_ID$'\t'"oasisia"$'\t'$annotations >> or
thogroups_annotations_originated_vestimentifera_cl1_oasisia.csv
33         done
34     else
35         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_originated_vestimentifera_
cl1_oasisia.csv

```

```
36     fi
37 done < orthogroups_gene_IDs_originated_vestimentifera_cl1_oasisi
a_OK.txt
38
39 while read line; do
40     genes=$(cut -f 2 <<< "$line")
41     echo $genes
42     orthogroup_ID=$(cut -f 1 <<< "$line")
43     echo $orthogroup_ID
44     if [[ "$genes" == RPAC* ]]
45     then
46         IFS=', '      # space is set as delimiter
47         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
48         for gene in "${ADDR[@]}; do
49             annotations=$(cut -f 13,14,15,18,20,24 ../../riftia_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
50             echo $orthogroup_ID$'\t'"riftia"$'\t'$annotations >> ort
hogroups_annotations_originated_vestimentifera_cl1_riftia.csv
51         done
52     else
53         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_originated_vestimentifera_
cl1_riftia.csv
54     fi
55 done < orthogroups_gene_IDs_originated_vestimentifera_cl1_riftia
_OK.txt
56
57 while read line; do
58     genes=$(cut -f 2 <<< "$line")
59     echo $genes
60     orthogroup_ID=$(cut -f 1 <<< "$line")
```

```

61     echo $orthogroup_ID
62     if [[ "$genes" == nbis* ]]
63     then
64         IFS=', '      # space is set as delimiter
65         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
66         for gene in "${ADDR[@]}; do
67             annotations=$(cut -f 13,14,15,18,20,24 ../../paraescarpia
_annotation_Jun2021_TrinoPantherK0.xls | fgrep $gene)
68             echo $orthogroup_ID$'\t'"paraescarpia"$'\t'$annotations
>> orthogroups_annotations_originated_vestimentifera_cl1_paraesc
arpia.csv
69         done
70     else
71         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_originated_vestimentifera_
cl1_paraescarpia.csv
72     fi
73 done < orthogroups_gene_IDs_originated_vestimentifera_cl1_paraes
carpia_OK.txt
74
75 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t
't'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
mber" > orthogroups_annotations_originated_vestimentifera_cl1_Oa
lv_Rpac_Pech.csv
76 cat orthogroups_annotations_originated_vestimentifera_cl1_oasisi
a.csv orthogroups_annotations_originated_vestimentifera_cl1_rift
ia.csv orthogroups_annotations_originated_vestimentifera_cl1_par
aescarpia.csv >> orthogroups_annotations_originated_vestimentife
ra_cl1_Oalv_Rpac_Pech.csv
77 sort orthogroups_annotations_originated_vestimentifera_cl1_Oalv_
Rpac_Pech.csv > orthogroups_annotations_originated_vestimentifer
a_cl1_Oalv_Rpac_Pech_OK.csv

```

```

78
79 echo "Orthogroup"${'\t'"Panther_annotation" > orthogroups_mostAbu
ndantAnnotation_originated_vestimentifera_cl1.csv
80 while read line; do
81     orthogroup_ID=$(cut -f 1 <<< "$line")
82     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_ori
ginated_vestimentifera_cl1_Ofra_Oalv_Rpac_Lluy_Pech_OK.csv | cut
-f 7 | sed '/^$/d' | sort | uniq -c | sort -r | awk '{s1=""; pr
int $0}' | head -1)
83     echo $orthogroup_ID${'\t'}$annotation >> orthogroups_mostA
bundantAnnotation_originated_vestimentifera_cl1.csv
84 done < gene_families_originated_vestimentifera_cl1.csv

```

orthogroups_annotations_originated_vestimentifera_cl2.sh

```

1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
ul2021/gains/vestimentifera_cl2
3  #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
l2021/gains/vestimentifera_cl2
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=100G
7  #$ -l h_rt=140:00:0
8  #$ -l highmem
9
10 cut -f 1,4 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
grep -w Vestimentifera_cl2 | cut -f 1 > gene_families_originated
_vestimentifera_cl2.txt #families originated in oasisia
11 fgrep -f gene_families_originated_vestimentifera_cl2.txt ../../O
rthogroups.csv > gene_families_originated_vestimentifera_cl2.csv
12
13 cut -f 1,19 gene_families_originated_vestimentifera_cl2.csv > or
thogroups_gene_IDs_originated_vestimentifera_cl2_oasisia.txt

```

```

14 sed 's/Oalv|//g' orthogroups_gene_IDs_originated_vestimentifera_
    cl2_oasisia.txt > orthogroups_gene_IDs_originated_vestimentifera_
    cl2_oasisia_OK.txt
15 cut -f 1,24 gene_families_originated_vestimentifera_cl2.csv > or
    thogroups_gene_IDs_originated_vestimentifera_cl2_riftia.txt
16 sed 's/Rpac|//g' orthogroups_gene_IDs_originated_vestimentifera_
    cl2_riftia.txt > orthogroups_gene_IDs_originated_vestimentifera_
    cl2_riftia_OK.txt
17
18
19
20 while read line; do
21     genes=$(cut -f 2 <<< "$line")
22     echo $genes
23     orthogroup_ID=$(cut -f 1 <<< "$line")
24     echo $orthogroup_ID
25     if [[ "$genes" == OALV* ]]
26     then
27         IFS=', '      # space is set as delimiter
28         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
29         for gene in "${ADDR[@]}; do
30             annotations=$(cut -f 13,14,15,18,20,24 ../../oasisia_anno
tation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
31             echo $orthogroup_ID$'\t'"oasisia"$'\t'$annotations >> or
thogroups_annotations_originated_vestimentifera_cl2_oasisia.csv
32         done
33     else
34         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_originated_vestimentifera_
cl2_oasisia.csv
35     fi

```



```

36 done < orthogroups_gene_IDs_originated_vestimentifera_cl2_oasisi
   a_OK.txt
37
38 while read line; do
39     genes=$(cut -f 2 <<< "$line")
40     echo $genes
41     orthogroup_ID=$(cut -f 1 <<< "$line")
42     echo $orthogroup_ID
43     if [[ "$genes" == RPAC* ]]
44     then
45         IFS=', '      # space is set as delimiter
46         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
47         for gene in "${ADDR[@]}"; do
48             annotations=$(cut -f 13,14,15,18,20,24 ../../riftia_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
49             echo $orthogroup_ID$'\t'"riftia"$'\t'$annotations >> ort
hogroups_annotations_originated_vestimentifera_cl2_riftia.csv
50         done
51     else
52         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_originated_vestimentifera_
cl2_riftia.csv
53     fi
54 done < orthogroups_gene_IDs_originated_vestimentifera_cl2_riftia
   _OK.txt
55
56 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t
't'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
mber" > orthogroups_annotations_originated_vestimentifera_cl2_Oa
lv_Rpac.csv
57 cat orthogroups_annotations_originated_vestimentifera_cl2_oasisi

```

```
a.csv orthogroups_annotations_originated_vestimentifera_cl2_rift
ia.csv >> orthogroups_annotations_originated_vestimentifera_cl2_
0alv_Rpac.csv
58 sort orthogroups_annotations_originated_vestimentifera_cl2_0alv_
Rpac.csv > orthogroups_annotations_originated_vestimentifera_cl2_
_0alv_Rpac_OK.csv
59
60 echo "Orthogroup"$(\t)"Panther_annotation" > orthogroups_mostAbu
ndantAnnotation_originated_vestimentifera_cl2.csv
61 while read line; do
62     orthogroup_ID=$(cut -f 1 <<< "$line")
63     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_ori
ginated_vestimentifera_cl2_0alv_Rpac_OK.csv | cut -f 7 | sed '/
^$/d' | sort | uniq -c | sort -r | awk '{s1=""; print $0}' | he
ad -1)
64     echo $orthogroup_ID$(\t)$annotation >> orthogroups_mostA
bundantAnnotation_originated_vestimentifera_cl2.csv
65 done < gene_families_originated_vestimentifera_cl2.csv
```

Losses

orthogroups_annotations_losses_osedax.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
ul2021/losses/osedax
3 #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
l2021/losses/osedax
4 #$ -j y
5 #$ -pe smp 1
6 #$ -l h_vmem=100G
7 #$ -l h_rt=140:00:0
8 #$ -l highmem
9
```

```
10 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |  
grep -w Ofra | cut -f 1 > gene_families_losses_osedax.txt #famil  
ies losses in oasisia  
11 fgrep -f gene_families_losses_osedax.txt ../../Orthogroups.csv >  
gene_families_losses_osedax.csv  
12  
13 cut -f 1,19 gene_families_losses_osedax.csv > orthogroups_gene_I  
Ds_losses_osedax_oasisia.txt  
14 sed 's/Oalv|//g' orthogroups_gene_IDs_losses_osedax_oasisia.txt  
> orthogroups_gene_IDs_losses_osedax_oasisia_OK.txt  
15 cut -f 1,24 gene_families_losses_osedax.csv > orthogroups_gene_I  
Ds_losses_osedax_riftia.txt  
16 sed 's/Rpac|//g' orthogroups_gene_IDs_losses_osedax_riftia.txt >  
orthogroups_gene_IDs_losses_osedax_riftia_OK.txt  
17 cut -f 1,14 gene_families_losses_osedax.csv > orthogroups_gene_I  
Ds_losses_osedax_lamellibrachia.txt  
18 sed 's/Lluy|//g' orthogroups_gene_IDs_losses_osedax_lamellibrach  
ia.txt > orthogroups_gene_IDs_losses_osedax_lamellibrachia_OK.tx  
t  
19 cut -f 1,23 gene_families_losses_osedax.csv > orthogroups_gene_I  
Ds_losses_osedax_paraescarpia.txt  
20 sed 's/Pech|//g' orthogroups_gene_IDs_losses_osedax_paraescarp  
ia.txt > orthogroups_gene_IDs_losses_osedax_paraescarpia_OK.txt  
21 cut -f 1,21 gene_families_losses_osedax.csv > orthogroups_gene_I  
Ds_losses_osedax_owenia.txt  
22 sed 's/Ofus|//g' orthogroups_gene_IDs_losses_osedax_owenia.txt >  
orthogroups_gene_IDs_losses_osedax_owenia_OK.txt  
23 cut -f 1,5 gene_families_losses_osedax.csv > orthogroups_gene_ID  
s_losses_osedax_capitella.txt  
24 sed 's/Ctel|//g' orthogroups_gene_IDs_losses_osedax_capitella.tx  
t > orthogroups_gene_IDs_losses_osedax_capitella_OK.txt  
25  
26
```

```

27 while read line; do
28     genes=$(cut -f 2 <<< "$line")
29     echo $genes
30     orthogroup_ID=$(cut -f 1 <<< "$line")
31     echo $orthogroup_ID
32     if [[ "$genes" == OFUS* ]]
33     then
34         IFS=', '      # space is set as delimiter
35         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
36         for gene in "${ADDR[@]}; do
37             cut -f 1,2,3,11,12,13 ../../Owenia_annotation_v250920.1_
TrinoPantherK0.xls | fgrep $gene > temp_file.txt
38             #K0_number=$(cut -f 1 temp_file.txt)
39             #gene_ID=$(cut -f 2 temp_file.txt)
40             #Panther_annotation=$(cut -f 3 temp_file.txt)
41             #GO_1=$(cut -f 4 temp_file.txt)
42             #GO_1=$(cut -f 5 temp_file.txt)
43             #GO_1=$(cut -f 6 temp_file.txt)
44             echo $orthogroup_ID$'\t'"owenia"$'\t'$(cut -f 4 temp_fil
e.txt)$'\t'$(cut -f 5 temp_file.txt)$'\t'$(cut -f 6 temp_file.tx
t)$'\t'$(cut -f 2 temp_file.txt)$'\t'$(cut -f 3 temp_file.tx
t)$'\t'$(cut -f 1 temp_file.txt) >> orthogroups_annotations_loss
es_osedax_owenia.csv
45         done
46     else
47         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_osedax_owenia.csv
48     fi
49 done < orthogroups_gene_IDs_losses_osedax_owenia_OK.txt
50
51 while read line; do

```

```

52     genes=$(cut -f 2 <<< "$line")
53     echo $genes
54     orthogroup_ID=$(cut -f 1 <<< "$line")
55     echo $orthogroup_ID
56     if [[ "$genes" == CapteT* ]]
57     then
58         IFS=', '      # space is set as delimiter
59         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
60         for gene in "${ADDR[@]}; do
61             cut -f 1,3,7,21,22 ../../Capitella_annotation_Feb2021_Tri
noPantherK0.xls | fgrep $gene > temp_file.txt
62             #K0_number=$(cut -f 7 temp_file.txt)
63             #gene_ID=$(cut -f 1 temp_file.txt)
64             #Panther_annotation=$(cut -f 3 temp_file.txt)
65             #GO_1=$(cut -f 21 temp_file.txt)
66             #GO_1=$(cut -f 22 temp_file.txt)
67             #GO_1=$(cut -f 6 temp_file.txt) NONE
68             echo $orthogroup_ID$'\t'"capitella"$'\t'$(cut -f 21 temp
_file.txt)$'\t'$(cut -f 22 temp_file.txt)$'\t'""'\t'$(cut -f 1 t
emp_file.txt)$'\t'$(cut -f 3 temp_file.txt)$'\t'$(cut -f 7 temp_
file.txt) >> orthogroups_annotations_losses_osedax_capitella.csv
69         done
70     else
71         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_osedax_capitella.cs
v
72     fi
73 done < orthogroups_gene_IDs_losses_osedax_capitella_0K.txt
74
75 while read line; do
76     genes=$(cut -f 2 <<< "$line")

```

```
77     echo $genes
78     orthogroup_ID=$(cut -f 1 <<< "$line")
79     echo $orthogroup_ID
80     if [[ "$genes" == OALV* ]]
81     then
82         IFS=', '      # space is set as delimiter
83         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
84         for gene in "${ADDR[@]}"; do
85             annotations=$(cut -f 13,14,15,18,20,24 ../../oasisia_anno
tation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
86             echo $orthogroup_ID$'\t'"oasisia"$'\t'$annotations >> or
thogroups_annotations_losses_osedax_oasisia.csv
87         done
88     else
89         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_osedax_oasisia.csv
90     fi
91 done < orthogroups_gene_IDs_losses_osedax_oasisia_OK.txt
92
93 while read line; do
94     genes=$(cut -f 2 <<< "$line")
95     echo $genes
96     orthogroup_ID=$(cut -f 1 <<< "$line")
97     echo $orthogroup_ID
98     if [[ "$genes" == RPAC* ]]
99     then
100         IFS=', '      # space is set as delimiter
101         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
102         for gene in "${ADDR[@]}"; do
```

```

103     annotations=$(cut -f 13,14,15,18,20,24 ../../riftia_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
104     echo $orthogroup_ID$'\t'"riftia"$'\t'$annotations >> ort
hogroups_annotations_losses_osedax_riftia.csv
105     done
106     else
107     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_osedax_riftia.csv
108     fi
109 done < orthogroups_gene_IDs_losses_osedax_riftia_OK.txt
110
111 while read line; do
112     genes=$(cut -f 2 <<< "$line")
113     echo $genes
114     orthogroup_ID=$(cut -f 1 <<< "$line")
115     echo $orthogroup_ID
116     if [[ "$genes" == FUN* ]]
117     then
118         IFS=', '      # space is set as delimiter
119         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
120         for gene in "${ADDR[@]}; do
121             annotations=$(cut -f 13,14,15,18,20,24 ../../lamellibrach
ia_annotation_Feb2021_TrinoPantherK0_OK.xls | fgrep $gene)
122             echo $orthogroup_ID$'\t'"lamellibrachia"$'\t'$annotation
s >> orthogroups_annotations_losses_osedax_lamellibrachia.csv
123         done
124     else
125     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_osedax_lamellibrach
ia.csv
126     fi

```

```

127 done < orthogroups_gene_IDs_losses_osedax_lamellibrachia_OK.txt
128
129 while read line; do
130     genes=$(cut -f 2 <<< "$line")
131     echo $genes
132     orthogroup_ID=$(cut -f 1 <<< "$line")
133     echo $orthogroup_ID
134     if [[ "$genes" == nbis* ]]
135     then
136         IFS=', '      # space is set as delimiter
137         read -ra ADDR <<< "$genes"  # str is read into an array a
s tokens separated by IFS
138         for gene in "${ADDR[@]}; do
139             annotations=$(cut -f 13,14,15,18,20,24 ../../paraescarpia
_annotation_Jun2021_TrinoPantherK0.xls | fgrep $gene)
140             echo $orthogroup_ID$'\t'"paraescarpia"$'\t'$annotations
>> orthogroups_annotations_losses_osedax_paraescarpia.csv
141         done
142     else
143         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
""$'\t'"" >> orthogroups_annotations_losses_osedax_paraescarpi
a.csv
144     fi
145 done < orthogroups_gene_IDs_losses_osedax_paraescarpia_OK.txt
146
147 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t'
t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
mber" > orthogroups_annotations_losses_osedax_Oalv_Rpac_Lluy_Pec
h_Ofus_Ctel.csv
148 cat orthogroups_annotations_losses_osedax_oasisia.csv orthogroup
s_annotations_losses_osedax_riftia.csv orthogroups_annotations_l
osses_osedax_lamellibrachia.csv orthogroups_annotations_losses_o

```



```

sedax_paraescarpia.csv orthogroups_annotations_losses_osedax_owe
nia.csv orthogroups_annotations_losses_osedax_capitella.csv >> o
rthogroups_annotations_losses_osedax_Oalv_Rpac_Lluy_Pech_Ofus_Ct
el.csv
149 sort orthogroups_annotations_losses_osedax_Oalv_Rpac_Lluy_Pech_0
fus_Ctel.csv > orthogroups_annotations_losses_osedax_Oalv_Rpac_L
luy_Pech_Ofus_Ctel_0K.csv
150
151 echo "Orthogroup"$('\t'"Panther_annotation" > orthogroups_mostAbu
ndantAnnotation_losses_osedax.csv
152 while read line; do
153     orthogroup_ID=$(cut -f 1 <<< "$line")
154     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_los
ses_osedax_Oalv_Rpac_Lluy_Pech_Ofus_Ctel_0K.csv | cut -f 7 | sed
'/^$/d' | sort | uniq -c | sort -r | awk '{ $1="" ; print $0 }' |
head -1)
155     echo $orthogroup_ID$('\t'$annotation >> orthogroups_mostA
bundantAnnotation_losses_osedax.csv
156 done < gene_families_losses_osedax.csv

```

orthogroups_annotations_losses_oasisia.sh

```

1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
ul2021/losses/oasisia
3  #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
l2021/losses/oasisia
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=100G
7  #$ -l h_rt=140:00:0
8  #$ -l highmem
9
10 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |

```

```
grep -w Oalv | cut -f 1 > gene_families_losses_oasisia.txt #fami
lies losses in oasisia
11 fgrep -f gene_families_losses_oasisia.txt ../../Orthogroups.csv
> gene_families_losses_oasisia.csv
12
13 cut -f 1,20 gene_families_losses_oasisia.csv > orthogroups_gene_
IDs_losses_oasisia_osedax.txt
14 sed 's/Ofra|//g' orthogroups_gene_IDs_losses_oasisia_osedax.txt
> orthogroups_gene_IDs_losses_oasisia_osedax_OK.txt
15 cut -f 1,24 gene_families_losses_oasisia.csv > orthogroups_gene_
IDs_losses_oasisia_riftia.txt
16 sed 's/Rpac|//g' orthogroups_gene_IDs_losses_oasisia_riftia.txt
> orthogroups_gene_IDs_losses_oasisia_riftia_OK.txt
17 cut -f 1,14 gene_families_losses_oasisia.csv > orthogroups_gene_
IDs_losses_oasisia_lamellibrachia.txt
18 sed 's/Lluy|//g' orthogroups_gene_IDs_losses_oasisia_lamellibrac
hia.txt > orthogroups_gene_IDs_losses_oasisia_lamellibrachia_OK.
txt
19 cut -f 1,23 gene_families_losses_oasisia.csv > orthogroups_gene_
IDs_losses_oasisia_paraescarpia.txt
20 sed 's/Pech|//g' orthogroups_gene_IDs_losses_oasisia_paraescarpia
a.txt > orthogroups_gene_IDs_losses_oasisia_paraescarpia_OK.txt
21 cut -f 1,21 gene_families_losses_oasisia.csv > orthogroups_gene_
IDs_losses_oasisia_owenia.txt
22 sed 's/Ofus|//g' orthogroups_gene_IDs_losses_oasisia_owenia.txt
> orthogroups_gene_IDs_losses_oasisia_owenia_OK.txt
23 cut -f 1,5 gene_families_losses_oasisia.csv > orthogroups_gene_I
Ds_losses_oasisia_capitella.txt
24 sed 's/Ctel|//g' orthogroups_gene_IDs_losses_oasisia_capitella.t
xt > orthogroups_gene_IDs_losses_oasisia_capitella_OK.txt
25
26
27 while read line; do
```

```
28 genes=$(cut -f 2 <<< "$line")
29 echo $genes
30 orthogroup_ID=$(cut -f 1 <<< "$line")
31 echo $orthogroup_ID
32 if [[ "$genes" == OFUS* ]]
33 then
34     IFS=', '      # space is set as delimiter
35     read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
36     for gene in "${ADDR[@]}; do
37         cut -f 1,2,3,11,12,13 ../Owenia_annotation_v250920.1_
TrinoPantherK0.xls | fgrep $gene > temp_file.txt
38         #K0_number=$(cut -f 1 temp_file.txt)
39         #gene_ID=$(cut -f 2 temp_file.txt)
40         #Panther_annotation=$(cut -f 3 temp_file.txt)
41         #GO_1=$(cut -f 4 temp_file.txt)
42         #GO_1=$(cut -f 5 temp_file.txt)
43         #GO_1=$(cut -f 6 temp_file.txt)
44         echo $orthogroup_ID$'\t'"owenia"$'\t'$(cut -f 4 temp_fil
e.txt)$'\t'$(cut -f 5 temp_file.txt)$'\t'$(cut -f 6 temp_file.tx
t)$'\t'$(cut -f 2 temp_file.txt)$'\t'$(cut -f 3 temp_file.tx
t)$'\t'$(cut -f 1 temp_file.txt) >> orthogroups_annotations_loss
es_oasisia_owenia.csv
45     done
46 else
47     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_oasisia_owenia.csv
48 fi
49 done < orthogroups_gene_IDs_losses_oasisia_owenia_OK.txt
50
51 while read line; do
52     genes=$(cut -f 2 <<< "$line")
```

```

53     echo $genes
54     orthogroup_ID=$(cut -f 1 <<< "$line")
55     echo $orthogroup_ID
56     if [[ "$genes" == CapteT* ]]
57     then
58         IFS=', '      # space is set as delimiter
59         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
60         for gene in "${ADDR[@]}"; do
61             cut -f 1,3,7,21,22 ../../Capitella_annotation_Feb2021_Tri
noPantherK0.xls | fgrep $gene > temp_file.txt
62             #K0_number=$(cut -f 7 temp_file.txt)
63             #gene_ID=$(cut -f 1 temp_file.txt)
64             #Panther_annotation=$(cut -f 3 temp_file.txt)
65             #GO_1=$(cut -f 21 temp_file.txt)
66             #GO_1=$(cut -f 22 temp_file.txt)
67             #GO_1=$(cut -f 6 temp_file.txt) NONE
68             echo $orthogroup_ID$'\t'"capitella"$'\t'$(cut -f 21 temp
_file.txt)$'\t'$(cut -f 22 temp_file.txt)$'\t'""$'\t'$(cut -f 1 t
emp_file.txt)$'\t'$(cut -f 3 temp_file.txt)$'\t'$(cut -f 7 temp_
file.txt) >> orthogroups_annotations_losses_oasisia_capitella.cs
v
69             done
70         else
71             echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_oasisia_capitella.c
sv
72         fi
73     done < orthogroups_gene_IDs_losses_oasisia_capitella_OK.txt
74
75     while read line; do
76         genes=$(cut -f 2 <<< "$line")

```

```
77     echo $genes
78     orthogroup_ID=$(cut -f 1 <<< "$line")
79     echo $orthogroup_ID
80     if [[ "$genes" == OFRA* ]]
81     then
82         IFS=', '      # space is set as delimiter
83         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
84         for gene in "${ADDR[@]}"; do
85             annotations=$(cut -f 13,14,15,18,20,24 ../../osedax_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
86             echo $orthogroup_ID$'\t'"osedax"$'\t'$annotations >> ort
hogroups_annotations_losses_oasisia_osedax.csv
87         done
88     else
89         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_oasisia_osedax.csv
90     fi
91 done < orthogroups_gene_IDs_losses_oasisia_osedax_OK.txt
92
93 while read line; do
94     genes=$(cut -f 2 <<< "$line")
95     echo $genes
96     orthogroup_ID=$(cut -f 1 <<< "$line")
97     echo $orthogroup_ID
98     if [[ "$genes" == RPAC* ]]
99     then
100         IFS=', '      # space is set as delimiter
101         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
102         for gene in "${ADDR[@]}"; do
```

```

103     annotations=$(cut -f 13,14,15,18,20,24 ../../riftia_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
104     echo $orthogroup_ID$'\t'"riftia"$'\t'$annotations >> ort
hogroups_annotations_losses_oasisia_riftia.csv
105     done
106     else
107     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_oasisia_riftia.csv
108     fi
109 done < orthogroups_gene_IDs_losses_oasisia_riftia_OK.txt
110
111 while read line; do
112     genes=$(cut -f 2 <<< "$line")
113     echo $genes
114     orthogroup_ID=$(cut -f 1 <<< "$line")
115     echo $orthogroup_ID
116     if [[ "$genes" == FUN* ]]
117     then
118         IFS=', '      # space is set as delimiter
119         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
120         for gene in "${ADDR[@]}; do
121             annotations=$(cut -f 13,14,15,18,20,24 ../../lamellibrach
ia_annotation_Feb2021_TrinoPantherK0_OK.xls | fgrep $gene)
122             echo $orthogroup_ID$'\t'"lamellibrachia"$'\t'$annotation
s >> orthogroups_annotations_losses_oasisia_lamellibrachia.csv
123         done
124     else
125     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_oasisia_lamellibrac
hia.csv
126     fi

```

```

127 done < orthogroups_gene_IDs_losses_oasisia_lamellibrachia_OK.txt
128
129 while read line; do
130     genes=$(cut -f 2 <<< "$line")
131     echo $genes
132     orthogroup_ID=$(cut -f 1 <<< "$line")
133     echo $orthogroup_ID
134     if [[ "$genes" == nbis* ]]
135     then
136         IFS=', '      # space is set as delimiter
137         read -ra ADDR <<< "$genes"  # str is read into an array a
s tokens separated by IFS
138         for gene in "${ADDR[@]}; do
139             annotations=$(cut -f 13,14,15,18,20,24 ../../paraescarpia
_annotation_Jun2021_TrinoPantherK0.xls | fgrep $gene)
140             echo $orthogroup_ID$'\t'"paraescarpia"$'\t'$annotations
>> orthogroups_annotations_losses_oasisia_paraescarpia.csv
141         done
142     else
143         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
""$'\t'"" >> orthogroups_annotations_losses_oasisia_paraescarpia
a.csv
144     fi
145 done < orthogroups_gene_IDs_losses_oasisia_paraescarpia_OK.txt
146
147 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t'
t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
mber" > orthogroups_annotations_losses_oasisia_Ofra_Rpac_Lluy_Pe
ch_Ofus_Ctel.csv
148 cat orthogroups_annotations_losses_oasisia_osedax.csv orthogroup
s_annotations_losses_oasisia_riftia.csv orthogroups_annotations_
losses_oasisia_lamellibrachia.csv orthogroups_annotations_losses

```

```

_oasisia_paraescarpia.csv orthogroups_annotations_losses_oasisia
_owenia.csv orthogroups_annotations_losses_oasisia_capitella.csv
>> orthogroups_annotations_losses_oasisia_Ofra_Rpac_Lluy_Pech_Of
us_Ctel.csv
149 sort orthogroups_annotations_losses_oasisia_Ofra_Rpac_Lluy_Pech_
Ofus_Ctel.csv > orthogroups_annotations_losses_oasisia_Ofra_Rpac
_Lluy_Pech_Ofus_Ctel_OK.csv
150
151 echo "Orthogroup"${'\t'"Panther_annotation" > orthogroups_mostAbu
ndantAnnotation_losses_oasisia.csv
152 while read line; do
153     orthogroup_ID=$(cut -f 1 <<< "$line")
154     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_los
ses_oasisia_Ofra_Rpac_Lluy_Pech_Ofus_Ctel_OK.csv | cut -f 7 | se
d '/^$/d' | sort | uniq -c | sort -r | awk '{ $1=""; print $0}'
| head -1)
155     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
bundantAnnotation_losses_oasisia.csv
156 done < gene_families_losses_oasisia.csv

```

orthogroups_annotations_losses_riftia.sh

```

1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
ul2021/losses/riftia
3  #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
l2021/losses/riftia
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=100G
7  #$ -l h_rt=140:00:0
8  #$ -l highmem
9
10 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |

```



```
grep -w Rpac | cut -f 1 > gene_families_losses_riftia.txt #families losses in oasisia
11 fgrep -f gene_families_losses_riftia.txt ../../Orthogroups.csv >
    gene_families_losses_riftia.csv
12
13 cut -f 1,20 gene_families_losses_riftia.csv > orthogroups_gene_IDS_losses_riftia_osedax.txt
14 sed 's/Ofra|//g' orthogroups_gene_IDS_losses_riftia_osedax.txt >
    orthogroups_gene_IDS_losses_riftia_osedax_OK.txt
15 cut -f 1,14 gene_families_losses_riftia.csv > orthogroups_gene_IDS_losses_riftia_lamellibrachia.txt
16 sed 's/Lluy|//g' orthogroups_gene_IDS_losses_riftia_lamellibrachia.txt >
    orthogroups_gene_IDS_losses_riftia_lamellibrachia_OK.txt
17 cut -f 1,23 gene_families_losses_riftia.csv > orthogroups_gene_IDS_losses_riftia_paraescarpia.txt
18 sed 's/Pech|//g' orthogroups_gene_IDS_losses_riftia_paraescarpia.txt >
    orthogroups_gene_IDS_losses_riftia_paraescarpia_OK.txt
19 cut -f 1,21 gene_families_losses_riftia.csv > orthogroups_gene_IDS_losses_riftia_owenia.txt
20 sed 's/Ofus|//g' orthogroups_gene_IDS_losses_riftia_owenia.txt >
    orthogroups_gene_IDS_losses_riftia_owenia_OK.txt
21 cut -f 1,5 gene_families_losses_riftia.csv > orthogroups_gene_IDS_losses_riftia_capitella.txt
22 sed 's/Ctel|//g' orthogroups_gene_IDS_losses_riftia_capitella.txt >
    orthogroups_gene_IDS_losses_riftia_capitella_OK.txt
23 cut -f 1,19 gene_families_losses_riftia.csv > orthogroups_gene_IDS_losses_riftia_oasisia.txt
24 sed 's/Oalv|//g' orthogroups_gene_IDS_losses_riftia_oasisia.txt >
    orthogroups_gene_IDS_losses_riftia_oasisia_OK.txt
25
26 while read line; do
27     genes=$(cut -f 2 <<< "$line")
```

```
28     echo $genes
29     orthogroup_ID=$(cut -f 1 <<< "$line")
30     echo $orthogroup_ID
31     if [[ "$genes" == OALV* ]]
32     then
33         IFS=', '      # space is set as delimiter
34         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
35         for gene in "${ADDR[@]}"; do
36             annotations=$(cut -f 13,14,15,18,20,24 ../../oasisia_anno
tation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
37             echo $orthogroup_ID$'\t'"oasisia"$'\t'$annotations >> or
thogroups_annotations_losses_riftia_oasisia.csv
38         done
39     else
40         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_riftia_oasisia.csv
41     fi
42 done < orthogroups_gene_IDs_losses_riftia_oasisia_OK.txt
43
44 while read line; do
45     genes=$(cut -f 2 <<< "$line")
46     echo $genes
47     orthogroup_ID=$(cut -f 1 <<< "$line")
48     echo $orthogroup_ID
49     if [[ "$genes" == OFUS* ]]
50     then
51         IFS=', '      # space is set as delimiter
52         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
53         for gene in "${ADDR[@]}"; do
```

```

54         cut -f 1,2,3,11,12,13 ../../Owenia_annotation_v250920.1_
TrinoPantherK0.xls | fgrep $gene > temp_file.txt
55         #K0_number=$(cut -f 1 temp_file.txt)
56         #gene_ID=$(cut -f 2 temp_file.txt)
57         #Panther_annotation=$(cut -f 3 temp_file.txt)
58         #GO_1=$(cut -f 4 temp_file.txt)
59         #GO_1=$(cut -f 5 temp_file.txt)
60         #GO_1=$(cut -f 6 temp_file.txt)
61         echo $orthogroup_ID$'\t'"owenia"$'\t'$(cut -f 4 temp_fil
e.txt)$'\t'$(cut -f 5 temp_file.txt)$'\t'$(cut -f 6 temp_file.tx
t)$'\t'$(cut -f 2 temp_file.txt)$'\t'$(cut -f 3 temp_file.tx
t)$'\t'$(cut -f 1 temp_file.txt) >> orthogroups_annotations_loss
es_riftia_owenia.csv
62         done
63     else
64         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_riftia_owenia.csv
65     fi
66 done < orthogroups_gene_IDs_losses_riftia_owenia_OK.txt
67
68 while read line; do
69     genes=$(cut -f 2 <<< "$line")
70     echo $genes
71     orthogroup_ID=$(cut -f 1 <<< "$line")
72     echo $orthogroup_ID
73     if [[ "$genes" == CapteT* ]]
74     then
75         IFS=', '      # space is set as delimiter
76         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
77         for gene in "${ADDR[@]}; do
78             cut -f 1,3,7,21,22 ../../Capitella_annotation_Feb2021_Tri

```

```

noPantherK0.xls | fgrep $gene > temp_file.txt
79     #K0_number=$(cut -f 7 temp_file.txt)
80     #gene_ID=$(cut -f 1 temp_file.txt)
81     #Panther_annotation=$(cut -f 3 temp_file.txt)
82     #GO_1=$(cut -f 21 temp_file.txt)
83     #GO_1=$(cut -f 22 temp_file.txt)
84     #GO_1=$(cut -f 6 temp_file.txt) NONE
85     echo $orthogroup_ID$'\t'"capitella"$'\t'$(cut -f 21 temp
_file.txt)$'\t'$(cut -f 22 temp_file.txt)$'\t'""'\t'$(cut -f 1 t
emp_file.txt)$'\t'$(cut -f 3 temp_file.txt)$'\t'$(cut -f 7 temp_
file.txt) >> orthogroups_annotations_losses_riftia_capitella.csv
86     done
87     else
88     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_riftia_capitella.cs
v
89     fi
90 done < orthogroups_gene_IDs_losses_riftia_capitella_0K.txt
91
92 while read line; do
93     genes=$(cut -f 2 <<< "$line")
94     echo $genes
95     orthogroup_ID=$(cut -f 1 <<< "$line")
96     echo $orthogroup_ID
97     if [[ "$genes" == OFRA* ]]
98     then
99         IFS=', '      # space is set as delimiter
100         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
101         for gene in "${ADDR[@]}; do
102             annotations=$(cut -f 13,14,15,18,20,24 ../../osedax_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)

```

```

103     echo $orthogroup_ID$'\t'"osedax"$'\t'$annotations >> ort
hogroups_annotations_losses_riftia_osedax.csv
104     done
105     else
106     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_riftia_osedax.csv
107     fi
108 done < orthogroups_gene_IDs_losses_riftia_osedax_OK.txt
109
110 while read line; do
111     genes=$(cut -f 2 <<< "$line")
112     echo $genes
113     orthogroup_ID=$(cut -f 1 <<< "$line")
114     echo $orthogroup_ID
115     if [[ "$genes" == FUN* ]]
116     then
117         IFS=', '      # space is set as delimiter
118         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
119         for gene in "${ADDR[@]}; do
120             annotations=$(cut -f 13,14,15,18,20,24 ../../lamellibrach
ia_annotation_Feb2021_TrinoPantherK0_OK.xls | fgrep $gene)
121             echo $orthogroup_ID$'\t'"lamellibrachia"$'\t'$annotation
s >> orthogroups_annotations_losses_riftia_lamellibrachia.csv
122         done
123     else
124     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_riftia_lamellibrach
ia.csv
125     fi
126 done < orthogroups_gene_IDs_losses_riftia_lamellibrachia_OK.txt
127

```

```

128 while read line; do
129     genes=$(cut -f 2 <<< "$line")
130     echo $genes
131     orthogroup_ID=$(cut -f 1 <<< "$line")
132     echo $orthogroup_ID
133     if [[ "$genes" == nbis* ]]
134     then
135         IFS=', '      # space is set as delimiter
136         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
137         for gene in "${ADDR[@]}; do
138             annotations=$(cut -f 13,14,15,18,20,24 ../../paraescarpia
_annotation_Jun2021_TrinoPantherK0.xls | fgrep $gene)
139             echo $orthogroup_ID$'\t'"paraescarpia"$'\t'$annotations
>> orthogroups_annotations_losses_riftia_paraescarpia.csv
140         done
141     else
142         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
""$'\t'"" >> orthogroups_annotations_losses_riftia_paraescarpi
a.csv
143     fi
144 done < orthogroups_gene_IDs_losses_riftia_paraescarpia_OK.txt
145
146 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t'
t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
mber" > orthogroups_annotations_losses_riftia_Ofra_Oalv_Lluy_Pec
h_Ofus_Ctel.csv
147 cat orthogroups_annotations_losses_riftia_oasisia.csv orthogroup
s_annotations_losses_riftia_osedax.csv orthogroups_annotations_l
osses_riftia_lamellibrachia.csv orthogroups_annotations_losses_r
iftia_paraescarpia.csv orthogroups_annotations_losses_riftia_owe
nia.csv orthogroups_annotations_losses_riftia_capitella.csv >> o
rthogroups_annotations_losses_riftia_Ofra_Oalv_Lluy_Pech_Ofus_Ct

```

```

el.csv
148 sort orthogroups_annotations_losses_riftia_0fra_0alv_Lluy_Pech_0
    fus_Ctel.csv > orthogroups_annotations_losses_riftia_0fra_0alv_L
    luy_Pech_0fus_Ctel_OK.csv
149
150 echo "Orthogroup"$'\t'"Panther_annotation" > orthogroups_mostAbu
    ndantAnnotation_losses_riftia.csv
151 while read line; do
152     orthogroup_ID=$(cut -f 1 <<< "$line")
153     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_lo
        ses_riftia_0fra_0alv_Lluy_Pech_0fus_Ctel_OK.csv | cut -f 7 | sed
        '/^$/d' | sort | uniq -c | sort -r | awk '{ $1="" ; print $0 }' |
        head -1)
154     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
        bundantAnnotation_losses_riftia.csv
155 done < gene_families_losses_riftia.csv

```

orthogroups_annotations_losses_lamellibrachia.sh

```

1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
    ul2021/losses/lamellibrachia
3  #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
    l2021/losses/lamellibrachia
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=100G
7  #$ -l h_rt=140:00:0
8  #$ -l highmem
9
10 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
    grep -w Lluy | cut -f 1 > gene_families_losses_lamellibrachia.tx
    t #families losses in oasisia
11 fgrep -f gene_families_losses_lamellibrachia.txt ../../Orthogrou

```

```
ps.csv > gene_families_losses_lamellibrachia.csv
12
13 cut -f 1,20 gene_families_losses_lamellibrachia.csv > orthogroup
s_gene_IDs_losses_lamellibrachia_osedax.txt
14 sed 's/Ofra|//g' orthogroups_gene_IDs_losses_lamellibrachia_osed
ax.txt > orthogroups_gene_IDs_losses_lamellibrachia_osedax_OK.tx
t
15 cut -f 1,23 gene_families_losses_lamellibrachia.csv > orthogroup
s_gene_IDs_losses_lamellibrachia_paraescarpia.txt
16 sed 's/Pech|//g' orthogroups_gene_IDs_losses_lamellibrachia_para
escarpia.txt > orthogroups_gene_IDs_losses_lamellibrachia_paraes
carpia_OK.txt
17 cut -f 1,21 gene_families_losses_lamellibrachia.csv > orthogroup
s_gene_IDs_losses_lamellibrachia_owenia.txt
18 sed 's/Ofus|//g' orthogroups_gene_IDs_losses_lamellibrachia_owen
ia.txt > orthogroups_gene_IDs_losses_lamellibrachia_owenia_OK.tx
t
19 cut -f 1,5 gene_families_losses_lamellibrachia.csv > orthogroups
_gene_IDs_losses_lamellibrachia_capitella.txt
20 sed 's/Ctel|//g' orthogroups_gene_IDs_losses_lamellibrachia_capi
tella.txt > orthogroups_gene_IDs_losses_lamellibrachia_capitella
_OK.txt
21 cut -f 1,19 gene_families_losses_lamellibrachia.csv > orthogroup
s_gene_IDs_losses_lamellibrachia_oasisia.txt
22 sed 's/Oalv|//g' orthogroups_gene_IDs_losses_lamellibrachia_oasi
sia.txt > orthogroups_gene_IDs_losses_lamellibrachia_oasisia_OK.
txt
23 cut -f 1,24 gene_families_losses_lamellibrachia.csv > orthogroup
s_gene_IDs_losses_lamellibrachia_riftia.txt
24 sed 's/Rpac|//g' orthogroups_gene_IDs_losses_lamellibrachia_rift
ia.txt > orthogroups_gene_IDs_losses_lamellibrachia_riftia_OK.tx
t
25
26 while read line; do
```



```

27     genes=$(cut -f 2 <<< "$line")
28     echo $genes
29     orthogroup_ID=$(cut -f 1 <<< "$line")
30     echo $orthogroup_ID
31     if [[ "$genes" == RPAC* ]]
32     then
33         IFS=', '      # space is set as delimiter
34         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
35         for gene in "${ADDR[@]}; do
36             annotations=$(cut -f 13,14,15,18,20,24 ../../riftia_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
37             echo $orthogroup_ID$'\t'"riftia"$'\t'$annotations >> ort
hogroups_annotations_losses_lamellibrachia_riftia.csv
38         done
39     else
40         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_lamellibrachia_rift
ia.csv
41     fi
42 done < orthogroups_gene_IDs_losses_lamellibrachia_riftia_OK.txt
43
44 while read line; do
45     genes=$(cut -f 2 <<< "$line")
46     echo $genes
47     orthogroup_ID=$(cut -f 1 <<< "$line")
48     echo $orthogroup_ID
49     if [[ "$genes" == OALV* ]]
50     then
51         IFS=', '      # space is set as delimiter
52         read -ra ADDR <<< "$genes"    # str is read into an array a

```

s tokens separated by IFS

```

53     for gene in "${ADDR[@]}"; do
54         annotations=$(cut -f 13,14,15,18,20,24 ../../oasisia_anno
tation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
55         echo $orthogroup_ID$'\t'"oasisia"$'\t'$annotations >> or
thogroups_annotations_losses_lamellibrachia_oasisia.csv
56     done
57     else
58         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_lamellibrachia_oasi
sia.csv
59     fi
60 done < orthogroups_gene_IDs_losses_lamellibrachia_oasisia_OK.txt
61
62 while read line; do
63     genes=$(cut -f 2 <<< "$line")
64     echo $genes
65     orthogroup_ID=$(cut -f 1 <<< "$line")
66     echo $orthogroup_ID
67     if [[ "$genes" == OFUS* ]]
68     then
69         IFS=', '      # space is set as delimiter
70         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
71         for gene in "${ADDR[@]}"; do
72             cut -f 1,2,3,11,12,13 ../../Owenia_annotation_v250920.1_
TrinoPantherK0.xls | fgrep $gene > temp_file.txt
73             #K0_number=$(cut -f 1 temp_file.txt)
74             #gene_ID=$(cut -f 2 temp_file.txt)
75             #Panther_annotation=$(cut -f 3 temp_file.txt)
76             #GO_1=$(cut -f 4 temp_file.txt)

```

```

77         #GO_1=$(cut -f 5 temp_file.txt)
78         #GO_1=$(cut -f 6 temp_file.txt)
79         echo $orthogroup_ID$'\t'"owenia"$'\t'$(cut -f 4 temp_file.txt)$'\t'$(cut -f 5 temp_file.txt)$'\t'$(cut -f 6 temp_file.txt)$'\t'$(cut -f 2 temp_file.txt)$'\t'$(cut -f 3 temp_file.txt)$'\t'$(cut -f 1 temp_file.txt) >> orthogroups_annotations_losses_lamellibrachia_owenia.csv
80     done
81     else
82         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'"" >> orthogroups_annotations_losses_lamellibrachia_owenia.csv
83     fi
84 done < orthogroups_gene_IDs_losses_lamellibrachia_owenia_0K.txt
85
86 while read line; do
87     genes=$(cut -f 2 <<< "$line")
88     echo $genes
89     orthogroup_ID=$(cut -f 1 <<< "$line")
90     echo $orthogroup_ID
91     if [[ "$genes" == CapteT* ]]
92     then
93         IFS=', '      # space is set as delimiter
94         read -ra ADDR <<< "$genes"    # str is read into an array as tokens separated by IFS
95         for gene in "${ADDR[@]}; do
96             cut -f 1,3,7,21,22 ../../Capitella_annotation_Feb2021_Tri
noPantherK0.xls | fgrep $gene > temp_file.txt
97             #K0_number=$(cut -f 7 temp_file.txt)
98             #gene_ID=$(cut -f 1 temp_file.txt)
99             #Panther_annotation=$(cut -f 3 temp_file.txt)
100            #GO_1=$(cut -f 21 temp_file.txt)

```

```

101     #GO_1=$(cut -f 22 temp_file.txt)
102     #GO_1=$(cut -f 6 temp_file.txt) NONE
103     echo $orthogroup_ID$'\t'"capitella"$'\t'$(cut -f 21 temp
_file.txt)$'\t'$(cut -f 22 temp_file.txt)$'\t'""'\t'$(cut -f 1 t
emp_file.txt)$'\t'$(cut -f 3 temp_file.txt)$'\t'$(cut -f 7 temp_
file.txt) >> orthogroups_annotations_losses_lamellibrachia_capit
ella.csv
104     done
105     else
106     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_lamellibrachia_capi
tella.csv
107     fi
108 done < orthogroups_gene_IDs_losses_lamellibrachia_capitella_OK.t
xt
109
110 while read line; do
111     genes=$(cut -f 2 <<< "$line")
112     echo $genes
113     orthogroup_ID=$(cut -f 1 <<< "$line")
114     echo $orthogroup_ID
115     if [[ "$genes" == OFRA* ]]
116     then
117         IFS=', '      # space is set as delimiter
118         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
119         for gene in "${ADDR[@]}"; do
120             annotations=$(cut -f 13,14,15,18,20,24 ../../osedax_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
121             echo $orthogroup_ID$'\t'"osedax"$'\t'$annotations >> ort
hogroups_annotations_losses_lamellibrachia_osedax.csv
122         done

```

```
123     else
124         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
        ""$'\t'"" >> orthogroups_annotations_losses_lamellibrachia_osed
        ax.csv
125     fi
126 done < orthogroups_gene_IDs_losses_lamellibrachia_osedax_0K.txt
127
128 while read line; do
129     genes=$(cut -f 2 <<< "$line")
130     echo $genes
131     orthogroup_ID=$(cut -f 1 <<< "$line")
132     echo $orthogroup_ID
133     if [[ "$genes" == nbis* ]]
134     then
135         IFS=', '      # space is set as delimiter
136         read -ra ADDR <<< "$genes"    # str is read into an array a
        s tokens separated by IFS
137         for gene in "${ADDR[@]}; do
138             annotations=$(cut -f 13,14,15,18,20,24 ../../paraescarpia
            _annotation_Jun2021_TrinoPantherK0.xls | fgrep $gene)
139             echo $orthogroup_ID$'\t'"paraescarpia"$'\t'$annotations
            >> orthogroups_annotations_losses_lamellibrachia_paraescarpia.cs
            v
140         done
141     else
142         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
        ""$'\t'"" >> orthogroups_annotations_losses_lamellibrachia_para
        escarpia.csv
143     fi
144 done < orthogroups_gene_IDs_losses_lamellibrachia_paraescarpia_0
        K.txt
145
```

```

146 echo "Orthogroup"'\t'"Species"'\t'"GO_term1"'\t'"GO_term2"'\t'"GO_term3"'\t'"gene_ID"'\t'"Panther_annotation"'\t'"KEGG_number" > orthogroups_annotations_losses_lamellibrachia_Ofra_Oalv_Rpac_Pech_Ofus_Ctel.csv
147 cat orthogroups_annotations_losses_lamellibrachia_oasisia.csv orthogroups_annotations_losses_lamellibrachia_osedax.csv orthogroups_annotations_losses_lamellibrachia_riftia.csv orthogroups_annotations_losses_lamellibrachia_paraescarpia.csv orthogroups_annotations_losses_lamellibrachia_owenia.csv orthogroups_annotations_losses_lamellibrachia_capitella.csv >> orthogroups_annotations_losses_lamellibrachia_Ofra_Oalv_Rpac_Pech_Ofus_Ctel.csv
148 sort orthogroups_annotations_losses_lamellibrachia_Ofra_Oalv_Rpac_Pech_Ofus_Ctel.csv > orthogroups_annotations_losses_lamellibrachia_Ofra_Oalv_Rpac_Pech_Ofus_Ctel_OK.csv
149
150 echo "Orthogroup"'\t'"Panther_annotation" > orthogroups_mostAbundantAnnotation_losses_lamellibrachia.csv
151 while read line; do
152     orthogroup_ID=$(cut -f 1 <<< "$line")
153     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_losses_lamellibrachia_Ofra_Oalv_Rpac_Pech_Ofus_Ctel_OK.csv | cut -f 7 | sed '/^$/d' | sort | uniq -c | sort -r | awk '{ $1="" }; print $0}' | head -1)
154     echo $orthogroup_ID'\t'$annotation >> orthogroups_mostAbundantAnnotation_losses_lamellibrachia.csv
155 done < gene_families_losses_lamellibrachia.csv

```

orthogroups_annotations_losses_paraescarpia.sh

```

1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/Jul2021/losses/paraescarpia
3 #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Jul2021/losses/paraescarpia
4 #$ -j y

```

```
5  #$ -pe smp 1
6  #$ -l h_vmem=100G
7  #$ -l h_rt=140:00:0
8  #$ -l highmem
9
10 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
    grep -w Pech | cut -f 1 > gene_families_losses_paraescarpia.txt
    #families losses in oasisia
11 fgrep -f gene_families_losses_paraescarpia.txt ../../Orthogroup
    s.csv > gene_families_losses_paraescarpia.csv
12
13 cut -f 1,20 gene_families_losses_paraescarpia.csv > orthogroups_
    gene_IDs_losses_paraescarpia_osedax.txt
14 sed 's/Ofra|//g' orthogroups_gene_IDs_losses_paraescarpia_oseda
    x.txt > orthogroups_gene_IDs_losses_paraescarpia_osedax_OK.txt
15 cut -f 1,21 gene_families_losses_paraescarpia.csv > orthogroups_
    gene_IDs_losses_paraescarpia_owenia.txt
16 sed 's/Ofus|//g' orthogroups_gene_IDs_losses_paraescarpia_oweni
    a.txt > orthogroups_gene_IDs_losses_paraescarpia_owenia_OK.txt
17 cut -f 1,5 gene_families_losses_paraescarpia.csv > orthogroups_g
    ene_IDs_losses_paraescarpia_capitella.txt
18 sed 's/Ctel|//g' orthogroups_gene_IDs_losses_paraescarpia_capite
    lla.txt > orthogroups_gene_IDs_losses_paraescarpia_capitella_OK.
    txt
19 cut -f 1,19 gene_families_losses_paraescarpia.csv > orthogroups_
    gene_IDs_losses_paraescarpia_oasisia.txt
20 sed 's/Oalv|//g' orthogroups_gene_IDs_losses_paraescarpia_oasisi
    a.txt > orthogroups_gene_IDs_losses_paraescarpia_oasisia_OK.txt
21 cut -f 1,24 gene_families_losses_paraescarpia.csv > orthogroups_
    gene_IDs_losses_paraescarpia_riftia.txt
22 sed 's/Rpac|//g' orthogroups_gene_IDs_losses_paraescarpia_rifti
    a.txt > orthogroups_gene_IDs_losses_paraescarpia_riftia_OK.txt
23 cut -f 1,14 gene_families_losses_paraescarpia.csv > orthogroups_
```

```

gene_IDs_losses_paraescarpia_lamellibrachia.txt
24 sed 's/Lluy|//g' orthogroups_gene_IDs_losses_paraescarpia_lamell
    ibrachia.txt > orthogroups_gene_IDs_losses_paraescarpia_lamellib
    rachia_OK.txt
25
26 while read line; do
27     genes=$(cut -f 2 <<< "$line")
28     echo $genes
29     orthogroup_ID=$(cut -f 1 <<< "$line")
30     echo $orthogroup_ID
31     if [[ "$genes" == FUN* ]]
32     then
33         IFS=', '      # space is set as delimiter
34         read -ra ADDR <<< "$genes"    # str is read into an array a
    s tokens separated by IFS
35         for gene in "${ADDR[@]}; do
36             annotations=$(cut -f 13,14,15,18,20,24 ../../lamellibrach
    ia_annotation_Feb2021_TrinoPantherKO_OK.xls | fgrep $gene)
37             echo $orthogroup_ID$'\t'"lamellibrachia"$'\t'$annotation
    s >> orthogroups_annotations_losses_paraescarpia_lamellibrachia.
    csv
38         done
39     else
40         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
    '""$'\t'"" >> orthogroups_annotations_losses_paraescarpia_lamell
    ibrachia.csv
41     fi
42 done < orthogroups_gene_IDs_losses_paraescarpia_lamellibrachia_O
    K.txt
43
44 while read line; do
45     genes=$(cut -f 2 <<< "$line")

```



```
46 echo $genes
47 orthogroup_ID=$(cut -f 1 <<< "$line")
48 echo $orthogroup_ID
49 if [[ "$genes" == RPAC* ]]
50 then
51     IFS=', '      # space is set as delimiter
52     read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
53     for gene in "${ADDR[@]}"; do
54         annotations=$(cut -f 13,14,15,18,20,24 ../../riftia_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
55         echo $orthogroup_ID$'\t'"riftia"$'\t'$annotations >> ort
hogroups_annotations_losses_paraescarpia_riftia.csv
56     done
57 else
58     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_paraescarpia_rifti
a.csv
59 fi
60 done < orthogroups_gene_IDs_losses_paraescarpia_riftia_OK.txt
61
62 while read line; do
63     genes=$(cut -f 2 <<< "$line")
64     echo $genes
65     orthogroup_ID=$(cut -f 1 <<< "$line")
66     echo $orthogroup_ID
67     if [[ "$genes" == OALV* ]]
68     then
69         IFS=', '      # space is set as delimiter
70         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
```

```

71     for gene in "${ADDR[@]}"; do
72         annotations=$(cut -f 13,14,15,18,20,24 ../../oasisia_anno
tation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
73         echo $orthogroup_ID$'\t'"oasisia"$'\t'$annotations >> or
thogroups_annotations_losses_paraescarpia_oasisia.csv
74     done
75     else
76         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_paraescarpia_oasisi
a.csv
77     fi
78 done < orthogroups_gene_IDs_losses_paraescarpia_oasisia_0K.txt
79
80 while read line; do
81     genes=$(cut -f 2 <<< "$line")
82     echo $genes
83     orthogroup_ID=$(cut -f 1 <<< "$line")
84     echo $orthogroup_ID
85     if [[ "$genes" == OFUS* ]]
86     then
87         IFS=', '      # space is set as delimiter
88         read -ra ADDR <<< "$genes"  # str is read into an array a
s tokens separated by IFS
89         for gene in "${ADDR[@]}"; do
90             cut -f 1,2,3,11,12,13 ../../Owenia_annotation_v250920.1_
TrinoPantherK0.xls | fgrep $gene > temp_file.txt
91             #K0_number=$(cut -f 1 temp_file.txt)
92             #gene_ID=$(cut -f 2 temp_file.txt)
93             #Panther_annotation=$(cut -f 3 temp_file.txt)
94             #GO_1=$(cut -f 4 temp_file.txt)
95             #GO_1=$(cut -f 5 temp_file.txt)

```

```
cut -f 1,3,7,21,22 ../../Capitella_annotation_Feb2021_Tri
noPantherK0.xls | fgrep $gene > temp_file.txt
```

```

120     #GO_1=$(cut -f 6 temp_file.txt) NONE
121     echo $orthogroup_ID$'\t'"capitella"$'\t'$(cut -f 21 temp
_file.txt)$'\t'$(cut -f 22 temp_file.txt)$'\t'""'\t'$(cut -f 1 t
emp_file.txt)$'\t'$(cut -f 3 temp_file.txt)$'\t'$(cut -f 7 temp_
file.txt) >> orthogroups_annotations_losses_paraescarpia_capitel
lla.csv
122     done
123     else
124     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_paraescarpia_capite
lla.csv
125     fi
126 done < orthogroups_gene_IDs_losses_paraescarpia_capitella_OK.txt
127
128 while read line; do
129     genes=$(cut -f 2 <<< "$line")
130     echo $genes
131     orthogroup_ID=$(cut -f 1 <<< "$line")
132     echo $orthogroup_ID
133     if [[ "$genes" == OFRA* ]]
134     then
135         IFS=', '      # space is set as delimiter
136         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
137         for gene in "${ADDR[@]}; do
138             annotations=$(cut -f 13,14,15,18,20,24 ../../osedax_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
139             echo $orthogroup_ID$'\t'"osedax"$'\t'$annotations >> ort
hogroups_annotations_losses_paraescarpia_osedax.csv
140         done
141     else
142     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t

```

```

142 '""$'\t'"" >> orthogroups_annotations_losses_paraescarpia_oseda
x.csv
143     fi
144 done < orthogroups_gene_IDs_losses_paraescarpia_osedax_OK.txt
145
146
147 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t'
t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
mber" > orthogroups_annotations_losses_paraescarpia_ofra_0alv_Rp
ac_Lluy_0fus_Ctel.csv
148 cat orthogroups_annotations_losses_paraescarpia_oasisia.csv orth
ogroups_annotations_losses_paraescarpia_osedax.csv orthogroups_a
nnotations_losses_paraescarpia_riftia.csv orthogroups_annotation
s_losses_paraescarpia_lamellibrachia.csv orthogroups_annotations
_losses_paraescarpia_owenia.csv orthogroups_annotations_losses_p
araescarpia_capitella.csv >> orthogroups_annotations_losses_para
escarpia_ofra_0alv_Rpac_Lluy_0fus_Ctel.csv
149 sort orthogroups_annotations_losses_paraescarpia_ofra_0alv_Rpac_
Lluy_0fus_Ctel.csv > orthogroups_annotations_losses_paraescarpia
_ofra_0alv_Rpac_Lluy_0fus_Ctel_OK.csv
150
151 echo "Orthogroup"$'\t'"Panther_annotation" > orthogroups_mostAbu
ndantAnnotation_losses_paraescarpia.csv
152 while read line; do
153     orthogroup_ID=$(cut -f 1 <<< "$line")
154     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_los
ses_paraescarpia_ofra_0alv_Rpac_Lluy_0fus_Ctel_OK.csv | cut -f 7
| sed '/^$/d' | sort | uniq -c | sort -r | awk '{ $1="" ; print
$0 }' | head -1)
155     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
bundantAnnotation_losses_paraescarpia.csv
156 done < gene_families_losses_paraescarpia.csv

```

orthogroups_annotations_losses_vestimentifera_cl2.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
   ul2021/losses/vestmentifera_cl2
3  #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
   l2021/losses/vestmentifera_cl2
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=100G
7  #$ -l h_rt=140:00:0
8  #$ -l highmem
9
10 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
   grep -w Vestimentifera_cl2 | cut -f 1 > gene_families_losses_ves
   timentifera_cl2.txt #families losses in oasisia
11 fgrep -f gene_families_losses_vestimentifera_cl2.txt ../../Ortho
   groups.csv > gene_families_losses_vestimentifera_cl2.csv
12
13 cut -f 1,20 gene_families_losses_vestimentifera_cl2.csv > orthog
   rroups_gene_IDs_losses_vestimentifera_cl2_osedax.txt
14 sed 's/Ofra|//g' orthogroups_gene_IDs_losses_vestimentifera_cl2_
   osedax.txt > orthogroups_gene_IDs_losses_vestimentifera_cl2_osed
   ax_OK.txt
15 cut -f 1,21 gene_families_losses_vestimentifera_cl2.csv > orthog
   rroups_gene_IDs_losses_vestimentifera_cl2_owenia.txt
16 sed 's/Ofus|//g' orthogroups_gene_IDs_losses_vestimentifera_cl2_
   owenia.txt > orthogroups_gene_IDs_losses_vestimentifera_cl2_owen
   ia_OK.txt
17 cut -f 1,5 gene_families_losses_vestimentifera_cl2.csv > orthogr
   oups_gene_IDs_losses_vestimentifera_cl2_capitella.txt
18 sed 's/Ctel|//g' orthogroups_gene_IDs_losses_vestimentifera_cl2_
   capitella.txt > orthogroups_gene_IDs_losses_vestimentifera_cl2_c
   apitella_OK.txt
19 cut -f 1,14 gene_families_losses_vestimentifera_cl2.csv > orthog
```

```

roups_gene_IDs_losses_vestimentifera_cl2_lamellibrachia.txt
20 sed 's/Lluy|//g' orthogroups_gene_IDs_losses_vestimentifera_cl2_
lamellibrachia.txt > orthogroups_gene_IDs_losses_vestimentifera_
cl2_lamellibrachia_OK.txt
21 cut -f 1,23 gene_families_losses_vestimentifera_cl2.csv > orthog
roups_gene_IDs_losses_vestimentifera_cl2_paraescarpia.txt
22 sed 's/Pech|//g' orthogroups_gene_IDs_losses_vestimentifera_cl2_
paraescarpia.txt > orthogroups_gene_IDs_losses_vestimentifera_cl
2_paraescarpia_OK.txt
23
24 while read line; do
25     genes=$(cut -f 2 <<< "$line")
26     echo $genes
27     orthogroup_ID=$(cut -f 1 <<< "$line")
28     echo $orthogroup_ID
29     if [[ "$genes" == nbis* ]]
30     then
31         IFS=', '      # space is set as delimiter
32         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
33         for gene in "${ADDR[@]}; do
34             annotations=$(cut -f 13,14,15,18,20,24 ../../paraescarpia
_annotation_Jun2021_TrinoPantherK0.xls | fgrep $gene)
35             echo $orthogroup_ID$'\t'"paraescarpia"$'\t'$annotations
>> orthogroups_annotations_losses_vestimentifera_cl2_paraescarp
a.csv
36         done
37     else
38         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_vestimentifera_cl2_
paraescarpia.csv
39     fi
40 done < orthogroups_gene_IDs_losses_vestimentifera_cl2_paraescarp

```

ia_OK.txt

```
41
42 while read line; do
43     genes=$(cut -f 2 <<< "$line")
44     echo $genes
45     orthogroup_ID=$(cut -f 1 <<< "$line")
46     echo $orthogroup_ID
47     if [[ "$genes" == FUN* ]]
48     then
49         IFS=', '      # space is set as delimiter
50         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
51         for gene in "${ADDR[@]}; do
52             annotations=$(cut -f 13,14,15,18,20,24 ../../lamellibrach
ia_annotation_Feb2021_TrinoPantherKO_OK.xls | fgrep $gene)
53             echo $orthogroup_ID$'\t'"lamellibrachia"$'\t'$annotation
s >> orthogroups_annotations_losses_vestimentifera_cl2_lamellibr
achia.csv
54         done
55     else
56         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_vestimentifera_cl2_
lamellibrachia.csv
57     fi
58 done < orthogroups_gene_IDs_losses_vestimentifera_cl2_lamellibra
chia_OK.txt
59
60 while read line; do
61     genes=$(cut -f 2 <<< "$line")
62     echo $genes
63     orthogroup_ID=$(cut -f 1 <<< "$line")
64     echo $orthogroup_ID
```



```

65     if [[ "$genes" == OFUS* ]]
66     then
67         IFS=', '          # space is set as delimiter
68         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
69         for gene in "${ADDR[@]}; do
70             cut -f 1,2,3,11,12,13 ../../Owenia_annotation_v250920.1_
TrinoPantherK0.xls | fgrep $gene > temp_file.txt
71             #K0_number=$(cut -f 1 temp_file.txt)
72             #gene_ID=$(cut -f 2 temp_file.txt)
73             #Panther_annotation=$(cut -f 3 temp_file.txt)
74             #GO_1=$(cut -f 4 temp_file.txt)
75             #GO_1=$(cut -f 5 temp_file.txt)
76             #GO_1=$(cut -f 6 temp_file.txt)
77             echo $orthogroup_ID$'\t'"owenia"$'\t'$(cut -f 4 temp_fil
e.txt)$'\t'$(cut -f 5 temp_file.txt)$'\t'$(cut -f 6 temp_file.tx
t)$'\t'$(cut -f 2 temp_file.txt)$'\t'$(cut -f 3 temp_file.tx
t)$'\t'$(cut -f 1 temp_file.txt) >> orthogroups_annotations_loss
es_vestimentifera_cl2_owenia.csv
78         done
79     else
80         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_vestimentifera_cl2_
owenia.csv
81     fi
82 done < orthogroups_gene_IDs_losses_vestimentifera_cl2_owenia_OK.
txt
83
84 while read line; do
85     genes=$(cut -f 2 <<< "$line")
86     echo $genes
87     orthogroup_ID=$(cut -f 1 <<< "$line")

```

```

88     echo $orthogroup_ID
89     if [[ "$genes" == CapteT* ]]
90     then
91         IFS=', '          # space is set as delimiter
92         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
93         for gene in "${ADDR[@]}; do
94             cut -f 1,3,7,21,22 ../../Capitella_annotation_Feb2021_Tri
noPantherK0.xls | fgrep $gene > temp_file.txt
95             #K0_number=$(cut -f 7 temp_file.txt)
96             #gene_ID=$(cut -f 1 temp_file.txt)
97             #Panther_annotation=$(cut -f 3 temp_file.txt)
98             #GO_1=$(cut -f 21 temp_file.txt)
99             #GO_1=$(cut -f 22 temp_file.txt)
100            #GO_1=$(cut -f 6 temp_file.txt) NONE
101            echo $orthogroup_ID$'\t'"capitella"$'\t'$(cut -f 21 temp
_file.txt)$'\t'$(cut -f 22 temp_file.txt)$'\t'""'\t'$(cut -f 1 t
emp_file.txt)$'\t'$(cut -f 3 temp_file.txt)$'\t'$(cut -f 7 temp_
file.txt) >> orthogroups_annotations_losses_vestimentifera_cl2_c
apitella.csv
102            done
103        else
104            echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_vestimentifera_cl2_
capitella.csv
105        fi
106    done < orthogroups_gene_IDs_losses_vestimentifera_cl2_capitella_
OK.txt
107
108    while read line; do
109        genes=$(cut -f 2 <<< "$line")
110        echo $genes

```

```

111 orthogroup_ID=$(cut -f 1 <<< "$line")
112 echo $orthogroup_ID
113 if [[ "$genes" == OFRA* ]]
114 then
115     IFS=', '      # space is set as delimiter
116     read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
117     for gene in "${ADDR[@]}"; do
118         annotations=$(cut -f 13,14,15,18,20,24 ../../osedax_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
119         echo $orthogroup_ID$'\t'"osedax"$'\t'$annotations >> ort
hogroups_annotations_losses_vestimentifera_cl2_osedax.csv
120     done
121 else
122     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_vestimentifera_cl2_
osedax.csv
123 fi
124 done < orthogroups_gene_IDs_losses_vestimentifera_cl2_osedax_OK.
txt
125
126
127 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t
't'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
mber" > orthogroups_annotations_losses_vestimentifera_cl2_0fra_P
ech_Lluy_Ofus_Ctel.csv
128 cat orthogroups_annotations_losses_vestimentifera_cl2_osedax.csv
orthogroups_annotations_losses_vestimentifera_cl2_paraescarpia.c
sv orthogroups_annotations_losses_vestimentifera_cl2_lamellibrac
hia.csv orthogroups_annotations_losses_vestimentifera_cl2_oweni
a.csv orthogroups_annotations_losses_vestimentifera_cl2_capitell
a.csv >> orthogroups_annotations_losses_vestimentifera_cl2_0fra_
Pech_Lluy_Ofus_Ctel.csv

```

```

129 sort orthogroups_annotations_losses_vestimentifera_cl2_0fra_Pech
    _Lluy_0fus_Ctel.csv > orthogroups_annotations_losses_vestimentif
    era_cl2_0fra_Pech_Lluy_0fus_Ctel_0K.csv
130
131 echo "Orthogroup"${'\t'"Panther_annotation" > orthogroups_mostAbu
    ndantAnnotation_losses_vestimentifera_cl2.csv
132 while read line; do
133     orthogroup_ID=$(cut -f 1 <<< "$line")
134     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_los
        ses_vestimentifera_cl2_0fra_Pech_Lluy_0fus_Ctel_0K.csv | cut -f
        7 | sed '/^$/d' | sort | uniq -c | sort -r | awk '{ $1="" ; print
        $0}' | head -1)
135     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
        bundantAnnotation_losses_vestimentifera_cl2.csv
136 done < gene_families_losses_vestimentifera_cl2.csv

```

orthogroups_annotations_losses_vestimentifera_cl1.sh

```

1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
    ul2021/losses/vestimentifera_cl1
3 #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
    l2021/losses/vestimentifera_cl1
4 #$ -j y
5 #$ -pe smp 1
6 #$ -l h_vmem=100G
7 #$ -l h_rt=140:00:0
8 #$ -l highmem
9
10 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
    grep -w Vestimentifera_cl1 | cut -f 1 > gene_families_losses_ves
    timentifera_cl1.txt #families losses in oasisia
11 fgrep -f gene_families_losses_vestimentifera_cl1.txt ../../Ortho
    groups.csv > gene_families_losses_vestimentifera_cl1.csv

```

```
12
13 cut -f 1,20 gene_families_losses_vestimentifera_cl1.csv > orthog
roups_gene_IDs_losses_vestimentifera_cl1_osedax.txt
14 sed 's/Ofra|//g' orthogroups_gene_IDs_losses_vestimentifera_cl1_
osedax.txt > orthogroups_gene_IDs_losses_vestimentifera_cl1_osed
ax_OK.txt
15 cut -f 1,21 gene_families_losses_vestimentifera_cl1.csv > orthog
roups_gene_IDs_losses_vestimentifera_cl1_owenia.txt
16 sed 's/Ofus|//g' orthogroups_gene_IDs_losses_vestimentifera_cl1_
owenia.txt > orthogroups_gene_IDs_losses_vestimentifera_cl1_owen
ia_OK.txt
17 cut -f 1,5 gene_families_losses_vestimentifera_cl1.csv > orthogr
roups_gene_IDs_losses_vestimentifera_cl1_capitella.txt
18 sed 's/Ctel|//g' orthogroups_gene_IDs_losses_vestimentifera_cl1_
capitella.txt > orthogroups_gene_IDs_losses_vestimentifera_cl1_c
apitella_OK.txt
19 cut -f 1,14 gene_families_losses_vestimentifera_cl1.csv > orthog
roups_gene_IDs_losses_vestimentifera_cl1_lamellibrachia.txt
20 sed 's/Lluy|//g' orthogroups_gene_IDs_losses_vestimentifera_cl1_
lamellibrachia.txt > orthogroups_gene_IDs_losses_vestimentifera_
cl1_lamellibrachia_OK.txt
21
22 while read line; do
23     genes=$(cut -f 2 <<< "$line")
24     echo $genes
25     orthogroup_ID=$(cut -f 1 <<< "$line")
26     echo $orthogroup_ID
27     if [[ "$genes" == FUN* ]]
28     then
29         IFS=', '      # space is set as delimiter
30         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
31         for gene in "${ADDR[@]}"; do
```

```

32     annotations=$(cut -f 13,14,15,18,20,24 ../../lamellibrach
ia_annotation_Feb2021_TrinoPantherK0_OK.xls | fgrep $gene)
33     echo $orthogroup_ID$'\t'"lamellibrachia"$'\t'$annotation
s >> orthogroups_annotations_losses_vestimentifera_cl1_lamellibr
achia.csv
34     done
35     else
36     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_vestimentifera_cl1_
lamellibrachia.csv
37     fi
38 done < orthogroups_gene_IDs_losses_vestimentifera_cl1_lamellibra
chia_OK.txt
39
40 while read line; do
41     genes=$(cut -f 2 <<< "$line")
42     echo $genes
43     orthogroup_ID=$(cut -f 1 <<< "$line")
44     echo $orthogroup_ID
45     if [[ "$genes" == OFUS* ]]
46     then
47         IFS=', '      # space is set as delimiter
48         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
49         for gene in "${ADDR[@]}; do
50             cut -f 1,2,3,11,12,13 ../../Owenia_annotation_v250920.1_
TrinoPantherK0.xls | fgrep $gene > temp_file.txt
51             #K0_number=$(cut -f 1 temp_file.txt)
52             #gene_ID=$(cut -f 2 temp_file.txt)
53             #Panther_annotation=$(cut -f 3 temp_file.txt)
54             #GO_1=$(cut -f 4 temp_file.txt)
55             #GO_1=$(cut -f 5 temp_file.txt)

```

```
#GO_1=$(cut -f 6 temp_file.txt)

echo $orthogroup_ID$'\t'"owenia"$'\t'$(cut -f 4 temp_file.txt)$'\t'$(cut -f 5 temp_file.txt)$'\t'$(cut -f 6 temp_file.txt)$'\t'$(cut -f 2 temp_file.txt)$'\t'$(cut -f 3 temp_file.txt)$'\t'$(cut -f 1 temp_file.txt) >> orthogroups_annotations_losses_vestimentifera_cll_owenia.csv

done

else

echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'"" >> orthogroups_annotations_losses_vestimentifera_cll_owenia.csv

fi

done < orthogroups_gene_IDS_losses_vestimentifera_cll_owenia_0K.txt

while read line; do

genes=$(cut -f 2 <<< "$line")

echo $genes

orthogroup_ID=$(cut -f 1 <<< "$line")

echo $orthogroup_ID

if [[ "$genes" == CapteT* ]]

then

IFS=', '      # space is set as delimiter

read -ra ADDR <<< "$genes"    # str is read into an array as tokens separated by IFS

for gene in "${ADDR[@]}; do

cut -f 1,3,7,21,22 ../../Capitella_annotation_Feb2021_Tri-noPantherK0.xls | fgrep $gene > temp_file.txt

#KO_number=$(cut -f 7 temp_file.txt)

#gene_ID=$(cut -f 1 temp_file.txt)

#Panther_annotation=$(cut -f 3 temp_file.txt)

#GO_1=$(cut -f 21 temp_file.txt)
```

```

79         #GO_1=$(cut -f 22 temp_file.txt)
80         #GO_1=$(cut -f 6 temp_file.txt) NONE
81         echo $orthogroup_ID$'\t'"capitella"$'\t'$(cut -f 21 temp
_file.txt)$'\t'$(cut -f 22 temp_file.txt)$'\t'""'\t'$(cut -f 1 t
emp_file.txt)$'\t'$(cut -f 3 temp_file.txt)$'\t'$(cut -f 7 temp_
file.txt) >> orthogroups_annotations_losses_vestimentifera_cl1_c
apitella.csv
82         done
83     else
84         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_vestimentifera_cl1_
capitella.csv
85     fi
86 done < orthogroups_gene_IDs_losses_vestimentifera_cl1_capitella_
OK.txt
87
88 while read line; do
89     genes=$(cut -f 2 <<< "$line")
90     echo $genes
91     orthogroup_ID=$(cut -f 1 <<< "$line")
92     echo $orthogroup_ID
93     if [[ "$genes" == OFRA* ]]
94     then
95         IFS=', '      # space is set as delimiter
96         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
97         for gene in "${ADDR[@]}; do
98             annotations=$(cut -f 13,14,15,18,20,24 ../../osedax_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
99             echo $orthogroup_ID$'\t'"osedax"$'\t'$annotations >> ort
hogroups_annotations_losses_vestimentifera_cl1_osedax.csv
100         done

```



```
101     else
102         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t'
        ""$'\t'"" >> orthogroups_annotations_losses_vestimentifera_cl1_
        osedax.csv
103     fi
104 done < orthogroups_gene_IDs_losses_vestimentifera_cl1_osedax_OK.
        txt
105
106
107 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t'
        t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
        mber" > orthogroups_annotations_losses_vestimentifera_cl1_ofra_L
        luy_ofus_Ctel.csv
108 cat orthogroups_annotations_losses_vestimentifera_cl1_osedax.csv
        orthogroups_annotations_losses_vestimentifera_cl1_lamellibrachi
        a.csv orthogroups_annotations_losses_vestimentifera_cl1_owenia.c
        sv orthogroups_annotations_losses_vestimentifera_cl1_capitella.c
        sv >> orthogroups_annotations_losses_vestimentifera_cl1_ofra_Lluy
        ofus_Ctel.csv
109 sort orthogroups_annotations_losses_vestimentifera_cl1_ofra_Lluy
        _ofus_Ctel.csv > orthogroups_annotations_losses_vestimentifera_c
        l1_ofra_Lluy_ofus_Ctel_OK.csv
110
111 echo "Orthogroup"$'\t'"Panther_annotation" > orthogroups_mostAbu
        ndantAnnotation_losses_vestimentifera_cl1.csv
112 while read line; do
113     orthogroup_ID=$(cut -f 1 <<< "$line")
114     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_lo
        ses_vestimentifera_cl1_ofra_Lluy_ofus_Ctel_OK.csv | cut -f 7 | s
        ed '/^$/d' | sort | uniq -c | sort -r | awk '{ $1="" ; print $0 }'
        | head -1)
115     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
        bundantAnnotation_losses_vestimentifera_cl1.csv
116 done < gene_families_losses_vestimentifera_cl1.csv
```

orthogroups_annotations_losses_vestimentifera.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
   ul2021/losses/vestimentifera
3  #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
   l2021/losses/vestimentifera
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=100G
7  #$ -l h_rt=140:00:0
8  #$ -l highmem
9
10 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
   grep -w Vestimentifera | cut -f 1 > gene_families_losses_vestime
   ntifera.txt #families losses in oasisia
11 fgrep -f gene_families_losses_vestimentifera.txt ../../Orthogrou
   ps.csv > gene_families_losses_vestimentifera.csv
12
13 cut -f 1,20 gene_families_losses_vestimentifera.csv > orthogroup
   s_gene_IDs_losses_vestimentifera_osedax.txt
14 sed 's/Ofra|//g' orthogroups_gene_IDs_losses_vestimentifera_osed
   ax.txt > orthogroups_gene_IDs_losses_vestimentifera_osedax_OK.tx
   t
15 cut -f 1,21 gene_families_losses_vestimentifera.csv > orthogroup
   s_gene_IDs_losses_vestimentifera_owenia.txt
16 sed 's/Ofus|//g' orthogroups_gene_IDs_losses_vestimentifera_owen
   ia.txt > orthogroups_gene_IDs_losses_vestimentifera_owenia_OK.tx
   t
17 cut -f 1,5 gene_families_losses_vestimentifera.csv > orthogroups
   _gene_IDs_losses_vestimentifera_capitella.txt
18 sed 's/Ctel|//g' orthogroups_gene_IDs_losses_vestimentifera_capi
   tella.txt > orthogroups_gene_IDs_losses_vestimentifera_capitella
```

_OK.txt

```

19
20 while read line; do
21     genes=$(cut -f 2 <<< "$line")
22     echo $genes
23     orthogroup_ID=$(cut -f 1 <<< "$line")
24     echo $orthogroup_ID
25     if [[ "$genes" == OFUS* ]]
26     then
27         IFS=', '      # space is set as delimiter
28         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
29         for gene in "${ADDR[@]}"; do
30             cut -f 1,2,3,11,12,13 ../../Owenia_annotation_v250920.1_
TrinoPantherK0.xls | fgrep $gene > temp_file.txt
31             #K0_number=$(cut -f 1 temp_file.txt)
32             #gene_ID=$(cut -f 2 temp_file.txt)
33             #Panther_annotation=$(cut -f 3 temp_file.txt)
34             #GO_1=$(cut -f 4 temp_file.txt)
35             #GO_1=$(cut -f 5 temp_file.txt)
36             #GO_1=$(cut -f 6 temp_file.txt)
37             echo $orthogroup_ID$'\t'"owenia"$'\t'$(cut -f 4 temp_fil
e.txt)$'\t'$(cut -f 5 temp_file.txt)$'\t'$(cut -f 6 temp_file.tx
t)$'\t'$(cut -f 2 temp_file.txt)$'\t'$(cut -f 3 temp_file.tx
t)$'\t'$(cut -f 1 temp_file.txt) >> orthogroups_annotations_loss
es_vestimentifera_owenia.csv
38         done
39     else
40         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_vestimentifera_owen
ia.csv
41     fi

```

```
42 done < orthogroups_gene_IDs_losses_vestimentifera_owenia_0K.txt
43
44 while read line; do
45     genes=$(cut -f 2 <<< "$line")
46     echo $genes
47     orthogroup_ID=$(cut -f 1 <<< "$line")
48     echo $orthogroup_ID
49     if [[ "$genes" == CapteT* ]]
50     then
51         IFS=', '          # space is set as delimiter
52         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
53         for gene in "${ADDR[@]}; do
54             cut -f 1,3,7,21,22 ../../Capitella_annotation_Feb2021_Tri
noPantherK0.xls | fgrep $gene > temp_file.txt
55             #K0_number=$(cut -f 7 temp_file.txt)
56             #gene_ID=$(cut -f 1 temp_file.txt)
57             #Panther_annotation=$(cut -f 3 temp_file.txt)
58             #GO_1=$(cut -f 21 temp_file.txt)
59             #GO_1=$(cut -f 22 temp_file.txt)
60             #GO_1=$(cut -f 6 temp_file.txt) NONE
61             echo $orthogroup_ID$'\t'"capitella"$'\t'$(cut -f 21 temp
_file.txt)$'\t'$(cut -f 22 temp_file.txt)$'\t'""'\t'$(cut -f 1 t
emp_file.txt)$'\t'$(cut -f 3 temp_file.txt)$'\t'$(cut -f 7 temp_
file.txt) >> orthogroups_annotations_losses_vestimentifera_capi
tella.csv
62         done
63     else
64         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_vestimentifera_capi
tella.csv
65     fi
```

```

66 done < orthogroups_gene_IDs_losses_vestimentifera_capitella_OK.t
    xt
67
68 while read line; do
69     genes=$(cut -f 2 <<< "$line")
70     echo $genes
71     orthogroup_ID=$(cut -f 1 <<< "$line")
72     echo $orthogroup_ID
73     if [[ "$genes" == OFRA* ]]
74     then
75         IFS=', '      # space is set as delimiter
76         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
77         for gene in "${ADDR[@]}; do
78             annotations=$(cut -f 13,14,15,18,20,24 ../../osedax_annot
ation_Jan2021_TrinoPantherK0.xls | fgrep $gene)
79             echo $orthogroup_ID$'\t'"osedax"$'\t'$annotations >> ort
hogroups_annotations_losses_vestimentifera_osedax.csv
80         done
81     else
82         echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
'""$'\t'"" >> orthogroups_annotations_losses_vestimentifera_osed
ax.csv
83     fi
84 done < orthogroups_gene_IDs_losses_vestimentifera_osedax_OK.txt
85
86
87 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\t
't'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
mber" > orthogroups_annotations_losses_vestimentifera_ofra_ofus_
Ctel.csv
88 cat orthogroups_annotations_losses_vestimentifera_osedax.csv ort

```

```
hogroups_annotations_losses_vestimentifera_owenia.csv orthogroup
s_annotations_losses_vestimentifera_capitella.csv >> orthogroups
_annotations_losses_vestimentifera_0fra_0fus_Ctel.csv
89 sort orthogroups_annotations_losses_vestimentifera_0fra_0fus_Cte
l.csv > orthogroups_annotations_losses_vestimentifera_0fra_0fus_
Ctel_0K.csv
90
91 echo "Orthogroup"$'\t'"Panther_annotation" > orthogroups_mostAbu
ndantAnnotation_losses_vestimentifera.csv
92 while read line; do
93     orthogroup_ID=$(cut -f 1 <<< "$line")
94     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_lo
ses_vestimentifera_0fra_0fus_Ctel_0K.csv | cut -f 7 | sed '/^$/d
' | sort | uniq -c | sort -r | awk '{ $1="" ; print $0 }' | head -
1)
95     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA
bundantAnnotation_losses_vestimentifera.csv
96 done < gene_families_losses_vestimentifera.csv
```

orthogroups_annotations_losses_siboglinidae.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/gene_family_evolution/ferdi_script/J
ul2021/losses/siboglinidae
3 #$ -o /data/scratch/btx654/gene_family_evolution/ferdi_script/Ju
l2021/losses/siboglinidae
4 #$ -j y
5 #$ -pe smp 1
6 #$ -l h_vmem=100G
7 #$ -l h_rt=140:00:0
8 #$ -l highmem
9
10 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
grep -w Siboglinidae | cut -f 1 > gene_families_losses_siboglini
```

```

dae.txt #families losses in oasisia
11 fgrep -f gene_families_losses_siboglinidae.txt ../../Orthogroup
s.csv > gene_families_losses_siboglinidae.csv
12
13 cut -f 1,21 gene_families_losses_siboglinidae.csv > orthogroups_
gene_IDs_losses_siboglinidae_owenia.txt
14 sed 's/Ofus|//g' orthogroups_gene_IDs_losses_siboglinidae_oweni
a.txt > orthogroups_gene_IDs_losses_siboglinidae_owenia_OK.txt
15 cut -f 1,5 gene_families_losses_siboglinidae.csv > orthogroups_g
ene_IDs_losses_siboglinidae_capitella.txt
16 sed 's/Ctel|//g' orthogroups_gene_IDs_losses_siboglinidae_capite
lla.txt > orthogroups_gene_IDs_losses_siboglinidae_capitella_OK.
txt
17
18 while read line; do
19     genes=$(cut -f 2 <<< "$line")
20     echo $genes
21     orthogroup_ID=$(cut -f 1 <<< "$line")
22     echo $orthogroup_ID
23     if [[ "$genes" == OFUS* ]]
24     then
25         IFS=', '      # space is set as delimiter
26         read -ra ADDR <<< "$genes"    # str is read into an array a
s tokens separated by IFS
27         for gene in "${ADDR[@]}; do
28             cut -f 1,2,3,11,12,13 ../../Owenia_annotation_v250920.1_
TrinoPantherK0.xls | fgrep $gene > temp_file.txt
29             #K0_number=$(cut -f 1 temp_file.txt)
30             #gene_ID=$(cut -f 2 temp_file.txt)
31             #Panther_annotation=$(cut -f 3 temp_file.txt)
32             #GO_1=$(cut -f 4 temp_file.txt)
33             #GO_1=$(cut -f 5 temp_file.txt)

```

```
cut -f 1,3,7,21,22 ../../Capitella_annotation_Feb2021_Tri
noPantherK0.xls | fgrep $gene > temp_file.txt
```



```

58     #GO_1=$(cut -f 6 temp_file.txt) NONE
59     echo $orthogroup_ID$'\t'"capitella"$'\t'$(cut -f 21 temp
    _file.txt)$'\t'$(cut -f 22 temp_file.txt)$'\t'""'\t'$(cut -f 1 t
    emp_file.txt)$'\t'$(cut -f 3 temp_file.txt)$'\t'$(cut -f 7 temp_
    file.txt) >> orthogroups_annotations_losses_siboglinidae_capitel
    la.csv
60     done
61     else
62     echo $orthogroup_ID$'\t'""$'\t'""$'\t'""$'\t'""$'\t'""$'\t
    '""$'\t'"" >> orthogroups_annotations_losses_siboglinidae_capite
    lla.csv
63     fi
64 done < orthogroups_gene_IDs_losses_siboglinidae_capitella_OK.txt
65
66
67 echo "Orthogroup"$'\t'"Species"$'\t'"GO_term1"$'\t'"GO_term2"$'\
    t'"GO_term3"$'\t'"gene_ID"$'\t'"Panther_annotation"$'\t'"KEGG_nu
    mber" > orthogroups_annotations_losses_siboglinidae_Ofus_Ctel.cs
    v
68 cat orthogroups_annotations_losses_siboglinidae_owenia.csv ortho
    groups_annotations_losses_siboglinidae_capitella.csv >> orthogro
    ups_annotations_losses_siboglinidae_Ofus_Ctel.csv
69 sort orthogroups_annotations_losses_siboglinidae_Ofus_Ctel.csv >
    orthogroups_annotations_losses_siboglinidae_Ofus_Ctel_OK.csv
70
71 echo "Orthogroup"$'\t'"Panther_annotation" > orthogroups_mostAbu
    ndantAnnotation_losses_siboglinidae.csv
72 while read line; do
73     orthogroup_ID=$(cut -f 1 <<< "$line")
74     annotation=$(fgrep $orthogroup_ID orthogroups_annotations_los
    ses_siboglinidae_Ofus_Ctel_OK.csv | cut -f 7 | sed '/^$/d' | sor
    t | uniq -c | sort -r | awk '{ $1="" ; print $0 }' | head -1)
75     echo $orthogroup_ID$'\t'$annotation >> orthogroups_mostA

```

```
bundantAnnotation_losses_siboglinidae.csv
```

```
76 done < gene_families_losses_siboglinidae.csv
```

GO_terms

Base code:

```
1 install.packages("BiocManager")
2 BiocManager::install("topGO")
3 install.packages("ggpubr")
4
5 library(topGO)
6 library(ggplot2)
7 library(ggpubr)
8 library(cowplot)
9
10 # Import gene universe: whole (GO-annotated) genome
11 geneID2GO <- readMappings(file = "/Users/giacomo/Desktop/R/GO_enrichment/osedax/osedax_GO_universe.txt") ### 21108 transcripts have GO annotation
12 geneUniverse <- names(geneID2GO)
13
14 # Import and transform genes of interest: 8 clusters from step2a
15 cluster1 <- read.table("/Users/giacomo/Desktop/R/GO_enrichment/osedax/gene_IDs_originated_osedax_Siboglinidae.txt", header=FALSE)
16 cluster1 <- as.character(cluster1$V1)
17 cluster1genelist <- factor(as.integer(geneUniverse %in% cluster1))
18 names(cluster1genelist) <- geneUniverse
19
20 # fisher testing of GO term enrichment for Molecular Function (MF)
```

```
21 #cluster 1 - red genes
22 cluster1_G0data_MF <- new("topG0data", description="Cluster1",
23                           ontology="MF", allGenes=cluster1geneli
24                           st,
25                           annot = annFUN.gene2G0, gene2G0 = gene
26                           ID2G0)
27 cluster1_resultFisher_MF <- runTest(cluster1_G0data_MF,
28                                     algorithm="classic", statist
29                                     ic="fisher")
30 cluster1_MF <- GenTable(cluster1_G0data_MF, classicFisher = clus
31                           ter1_resultFisher_MF,
32                           orderBy = "resultFisher", ranksOf = "cla
33                           ssicFisher", topNodes = 15)
34
35 cluster1_MF[cluster1_MF == "< 1e-30"] <- "1e-30"
36 cluster1_MF[cluster1_MF == "<1e-30"] <- "1e-30"
37
38 goEnrichment <- cluster1_MF
39 goEnrichment$classicFisher <- as.numeric(goEnrichment$classicFis
40                           her)
41 goEnrichment <- goEnrichment[,c("GO.ID", "Term", "classicFisher")]
42 goEnrichment$Term <- gsub(" [a-z]*\\.\\.\\.\\.$", "", goEnrichmen
43                           t$Term)
44 goEnrichment$Term <- gsub("\\\\.\\.\\.\\.$", "", goEnrichment$Term)
45 goEnrichment$Term <- paste(goEnrichment$GO.ID, goEnrichment$Ter
46                           m, sep=", ")
47 goEnrichment$Term <- factor(goEnrichment$Term, levels=rev(goEnri
48                           chment$Term))
49
50 #it could happen that the second line of the previous block will
51 #give the error " Warning message:NAs introduced by coercion "
52 # a fix for that is " goEnrichment$classicFisher <- c(30, 30, 3
```

```
0, 30, 30, 30, 30, 30, 30, 30, 30, 30, 30, 30, 30, 30) "
```

```
43
44 cluster_1_plot <- ggplot(goEnrichment, aes(x=Term, y=-log10(classicFisher))) +
45   stat_summary(geom = "bar", fun = mean, position = "dodge") +
46   xlab("Molecular Function") +
47   ylab("-log10(p-value)") +
48   ggtitle("GF gains in Siboglinidae (Osedax)") +
49   scale_y_continuous(limits=c(0,30),breaks=round(seq(0,30, by =
50 2), 1)) +
51   theme_classic() +
52   theme(
53     legend.position='none',
54     legend.background=element_rect(),
55     plot.title=element_text(angle=0, size=12, face="bold", vjust
56 =1),
57     axis.text.x=element_text(angle=0, size=10, hjust=1.10),
58     axis.text.y=element_text(angle=0, size=10, vjust=0.5),
59     axis.title=element_text(size=12),
60     legend.key=element_blank(),      #removes the border
61     legend.key.size=unit(1, "cm"),    #Sets overall area/size
62 of the legend
63     legend.text=element_text(size=12), #Text size
64     title=element_text(size=12)) +
65   guides(colour=guide_legend(override.aes=list(size=2.5))) +
66   coord_flip()
67
68 cluster_1_plot
69
70 cluster1_G0data_BP <- new("topG0data", description="Cluster1",
71   ontology="BP", allGenes=cluster1geneli
```

```

st,
69         annot = annFUN.gene2GO, gene2GO = gene
ID2GO)
70 cluster1_resultFisher_BP <- runTest(cluster1_G0data_BP,
71         algorithm="classic", statist
ic="fisher")
72 cluster1_BP <- GenTable(cluster1_G0data_BP, classicFisher = clus
ter1_resultFisher_BP,
73         orderBy = "resultFisher", ranksOf = "cla
ssicFisher", topNodes = 15)
74
75 cluster1_BP[cluster1_BP == "< 1e-30"] <- "1e-30"
76 cluster1_BP[cluster1_BP == "<1e-30"] <- "1e-30"
77
78 goEnrichment <- cluster1_BP
79 goEnrichment$classicFisher <- as.numeric(goEnrichment$classicFis
her)
80 goEnrichment <- goEnrichment[,c("GO.ID","Term","classicFisher")]
81 goEnrichment$Term <- gsub(" [a-z]*\\.\\.\\.\\.\\$", "", goEnrichmen
t$Term)
82 goEnrichment$Term <- gsub("\\\\.\\.\\.\\.\\$", "", goEnrichment$Term)
83 goEnrichment$Term <- paste(goEnrichment$GO.ID, goEnrichment$Ter
m, sep=", ")
84 goEnrichment$Term <- factor(goEnrichment$Term, levels=rev(goEnri
chment$Term))
85
86 #it could happen that the second line of the previous block will
give the error " Warning message:NAs introduced by coercion "
87 # a fix for that is " goEnrichment$classicFisher <- c(30, 30, 3
0, 30, 30, 30, 30, 30, 30, 30, 30, 30, 30, 30, 30) "
88

```

```

89 cluster_1BP_plot <- ggplot(goEnrichment, aes(x=Term, y=-log10(classicFisher))) +
90   stat_summary(geom = "bar", fun = mean, position = "dodge") +
91   xlab("Biological Process") +
92   ylab("-log10(p-value)") +
93   scale_y_continuous(limits=c(0,30),breaks=round(seq(0,30, by =
94     2), 1)) +
95   theme_classic() +
96   theme(
97     legend.position='none',
98     legend.background=element_rect(),
99     plot.title=element_text(angle=0, size=12, face="bold", vjust
100   =1),
101     axis.text.x=element_text(angle=0, size=10, hjust=1.10),
102     axis.text.y=element_text(angle=0, size=10, vjust=0.5),
103     axis.title=element_text(size=12),
104     legend.key=element_blank(),      #removes the border
105     legend.key.size=unit(1, "cm"),    #Sets overall area/size
of the legend
106     legend.text=element_text(size=12), #Text size
107     title=element_text(size=12)) +
108   guides(colour=guide_legend(override.aes=list(size=2.5))) +
109   coord_flip()
110
111 cluster_1BP_plot
112
113 plot_grid(cluster_1BP_plot + rremove("x.title"),
114           cluster_1BP_plot,
115           ncol = 1, align="v")

```

- exported in pdf 7x8 inches

Riftia example

```
1 cd /data/scratch/btx654/gene_family_evolution/ferdi_script/Jul20
  21/GO_terms/riftia
2 # I will use BlastX GO terms
3 cut -f2,13 ../../riftia_annotation_Jan2021_TrinoPantherK0.xls |
  tail -n +2 > riftia_GO_raw.txt # 38179 genes in riftia_GO_raw.tx
  t
4 grep 'GO' riftia_GO_raw.txt > riftia_GO_only_raw.txt # 20737 rif
  tia_GO_only_raw.txt
5
```

python.py

```
1 if __name__ == "__main__":
2
3     import re
4
5     i = open("riftia_GO_only_raw.txt", "r")
6     o = open("riftia_GO_universe.txt", "w")
7
8     regex = re.compile(r'GO:\d+')
9
10    for line in i:
11        GOMatches = regex.findall(line)
12        Gene_ID = line.split("\t",1)[0]
13        if not GOMatches == []:
14            o.write(Gene_ID+"\t")
15            for i, match in enumerate(GOMatches):
16                if i+1 == len(GOMatches):
17                    o.write(match.strip("'")+"\n")
18                else:
19                    o.write(match.strip("'")+", ")
```

```
1 module load python
2 python python.py
```

Now we have the file riftia_GO_universe.txt and we need to select a subgroup of genes:

Expansions

```
1 cut -f 1,7 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
  grep -w Rpac | cut -f 1 > gene_families_expanded_riftia.txt #fam
  ilies expanded in riftia
2 fgrep -f gene_families_expanded_riftia.txt ../../Orthogroups.csv
  > gene_families_expanded_riftia.csv
3 cut -f 1,24 gene_families_expanded_riftia.csv | sed 's/Rpac|//g'
  | cut -f 2 | sed 's/, /\n/g' > gene_IDs_expanded_riftia.txt
```

R script:

```
1 install.packages("BiocManager")
2 BiocManager::install("topGO")
3 install.packages("ggpubr")
4
5 library(topGO)
6 library(ggplot2)
7 library(ggpubr)
8 library(cowplot)
9
10 # Import gene universe: whole (GO-annotated) genome
11 geneID2GO <- readMappings(file = "/Users/giacomo/Desktop/R/GO_en
  richment/riftia_GO_universe.txt") ### 21108 transcripts have GO
  annotation
12 geneUniverse <- names(geneID2GO)
13
14 # Import and transform genes of interest: 8 clusters from step2a
15 cluster1 <- read.table("/Users/giacomo/Desktop/R/GO_enrichment/g
```



```

ene_IDs_expanded_riftia.txt",header=FALSE)
16 cluster1 <- as.character(cluster1$V1)
17 cluster1genelist <- factor(as.integer(geneUniverse %in% cluster
18 1))
19 names(cluster1genelist) <- geneUniverse
20 # fisher testing of GO term enrichment for Molecular Function (M
21 F)
22 #cluster 1 - red genes
23 cluster1_GOdata_MF <- new("topGOdata", description="Cluster1",
24 ontology="MF", allGenes=cluster1geneli
25 st,
26 annot = annFUN.gene2GO, gene2GO = gene
27 ID2GO)
28 cluster1_resultFisher_MF <- runTest(cluster1_GOdata_MF,
29 algorithm="classic", statist
30 ic="fisher")
31 cluster1_MF <- GenTable(cluster1_GOdata_MF, classicFisher = clus
32 ter1_resultFisher_MF,
33 orderBy = "resultFisher", ranksOf = "cla
34 ssicFisher", topNodes = 15)
35 cluster1_MF[cluster1_MF == "< 1e-30"] <- "1e-30"
36
37 goEnrichment <- cluster1_MF
38 goEnrichment$classicFisher <- as.numeric(goEnrichment$classicFis
39 her)
40 goEnrichment <- goEnrichment[,c("GO.ID","Term","classicFisher")]
41 goEnrichment$Term <- gsub(" [a-z]*\\.\\.\\.\\.$", "", goEnrichmen
42 t$Term)
43 goEnrichment$Term <- gsub("\\\\.\\.\\.\\.$", "", goEnrichment$Term)
44 goEnrichment$Term <- paste(goEnrichment$GO.ID, goEnrichment$Ter

```

```

m, sep=", ")
38 goEnrichment$Term <- factor(goEnrichment$Term, levels=rev(goEnri
chment$Term))
39
40
41 cluster_1_plot <- ggplot(goEnrichment, aes(x=Term, y=-log10(clas
sicFisher))) +
42   stat_summary(geom = "bar", fun = mean, position = "dodge") +
43   xlab("Biological process") +
44   ylab("-log10(p-value)") +
45   ggtitle("GF expansions in Riftia") +
46   scale_y_continuous(limits=c(0,30),breaks=round(seq(0,30, by =
2), 1)) +
47   theme_classic() +
48   theme(
49     legend.position='none',
50     legend.background=element_rect(),
51     plot.title=element_text(angle=0, size=12, face="bold", vjust
=1),
52     axis.text.x=element_text(angle=0, size=10, hjust=1.10),
53     axis.text.y=element_text(angle=0, size=10, vjust=0.5),
54     axis.title=element_text(size=12),
55     legend.key=element_blank(),      #removes the border
56     legend.key.size=unit(1, "cm"),    #Sets overall area/size
of the legend
57     legend.text=element_text(size=12), #Text size
58     title=element_text(size=12)) +
59     guides(colour=guide_legend(override.aes=list(size=2.5))) +
60     coord_flip()
61
62 cluster_1_plot

```

- exported in pdf 4x8 inches

Gains

```

1 cut -f 1,4 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
  grep -w Rpac | cut -f 1 > gene_families_originated_riftia.txt #f
  amilies expanded in riftia
2 fgrep -f gene_families_originated_riftia.txt ../../Orthogroups.c
  sv > gene_families_originated_riftia.csv
3 cut -f 1,24 gene_families_originated_riftia.csv | sed 's/Rpac|//
  g' | cut -f 2 | sed 's/, /\n/g' > gene_IDs_originated_riftia.txt
4
5 cut -f 1,4 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
  grep -w Siboglinidae | cut -f 1 > gene_families_originated_rifti
  a_siboglinidae.txt #families expanded in riftia
6 fgrep -f gene_families_originated_riftia_siboglinidae.txt ../../
  Orthogroups.csv > gene_families_originated_riftia_siboglinidae.c
  sv
7 cut -f 1,24 gene_families_originated_riftia_siboglinidae.csv | s
  ed 's/Rpac|//g' | cut -f 2 | sed 's/, /\n/g' | sed '/^$/d' > gen
  e_IDs_originated_riftia_siboglinidae.txt
8
9 cut -f 1,4 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
  grep -w Vestimentifera | cut -f 1 > gene_families_originated_rif
  tia_Vestimentifera.txt #families expanded in riftia
10 fgrep -f gene_families_originated_riftia_Vestimentifera.txt
  ../../Orthogroups.csv > gene_families_originated_riftia_Vestimen
  tifera.csv
11 cut -f 1,24 gene_families_originated_riftia_Vestimentifera.csv |
  sed 's/Rpac|//g' | cut -f 2 | sed 's/, /\n/g' | sed '/^$/d' > ge
  ne_IDs_originated_riftia_Vestimentifera.txt
12
13 cut -f 1,4 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
  grep -w Vestimentifera_cl1 | cut -f 1 > gene_families_originated
  _riftia_Vestimentifera_cl1.txt #families expanded in riftia

```

```

14 fgrep -f gene_families_originated_riftia_Vestimentifera_cl1.txt
   ../../Orthogroups.csv > gene_families_originated_riftia_Vestimen
   tifera_cl1.csv
15 cut -f 1,24 gene_families_originated_riftia_Vestimentifera_cl1.c
   sv | sed 's/Rpac|//g' | cut -f 2 | sed 's/, /\n/g' | sed '/^$/d'
   > gene_IDs_originated_riftia_Vestimentifera_cl1.txt
16
17 cut -f 1,4 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
   grep -w Vestimentifera_cl2 | cut -f 1 > gene_families_originated
   _riftia_Vestimentifera_cl2.txt #families expanded in riftia
18 fgrep -f gene_families_originated_riftia_Vestimentifera_cl2.txt
   ../../Orthogroups.csv > gene_families_originated_riftia_Vestimen
   tifera_cl2.csv
19 cut -f 1,24 gene_families_originated_riftia_Vestimentifera_cl2.c
   sv | sed 's/Rpac|//g' | cut -f 2 | sed 's/, /\n/g' | sed '/^$/d'
   > gene_IDs_originated_riftia_Vestimentifera_cl2.txt

```

Losses

using Owenia to annotate the losses

```

1 cd /data/scratch/btx654/gene_family_evolution/ferdi_script/Jul20
   21/GO_terms/losses
2 cut -f2,11 ../../Owenia_annotation_v250920.1_TrinoPantherK0.xls
   | tail -n +2 | sed 's/"//g' > owenia_GO_raw.txt #
3 grep 'GO' owenia_GO_raw.txt > owenia_GO_only_raw.txt # 21108 ose
   dax_GO_only_raw.txt

```

python.py

```

1 if __name__ == "__main__":
2
3     import re
4
5     i = open("owenia_GO_only_raw.txt", "r")
6     o = open("owenia_GO_universe.txt", "w")

```

```

7
8     regex = re.compile(r'GO:\d+')
9
10    for line in i:
11        GOMatches = regex.findall(line)
12        Gene_ID = line.split("\t",1)[0]
13        if not GOMatches == []:
14            o.write(Gene_ID+"\t")
15            for i, match in enumerate(GOMatches):
16                if i+1 == len(GOMatches):
17                    o.write(match.strip("'")+"\n")
18                else:
19                    o.write(match.strip("'")+", ")

```

```

1 module load python
2 python python.py

```

losses.sh

```

1 cut -f 1,21 ../../Orthogroups.csv | tail -n +2 | grep -w Ofus |
  cut -f 1 > owenia_all_GF
2
3 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
  grep -w Ofra | cut -f 1 > gene_families_lost_osedax.txt
4 grep -f gene_families_lost_osedax.txt owenia_all_GF > losses_ose
  dax_owenia
5 fgrep -f losses_osedax_owenia ../../Orthogroups.csv > gene_famil
  ies_lost_osedax_owenia.csv
6 cut -f 1,21 gene_families_lost_osedax_owenia.csv | sed 's/Ofu
  s|//g' | cut -f 2 | sed 's/, /\n/g' > gene_IDs_lost_osedax_oweni
  a.txt
7

```

```
8 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
grep -w Oalv | cut -f 1 > gene_families_lost_oasisia.txt
9 grep -f gene_families_lost_oasisia.txt owenia_all_GF > losses_oa
sisia_owenia
10 fgrep -f losses_oasisia_owenia ../../Orthogroups.csv > gene_fami
lies_lost_oasisia_owenia.csv
11 cut -f 1,21 gene_families_lost_oasisia_owenia.csv | sed 's/Ofu
s|//g' | cut -f 2 | sed 's/, /\n/g' > gene_IDs_lost_oasisia_owen
ia.txt
12
13 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
grep -w Rpac | cut -f 1 > gene_families_lost_riftia.txt
14 grep -f gene_families_lost_riftia.txt owenia_all_GF > losses_rif
tia_owenia
15 fgrep -f losses_riftia_owenia ../../Orthogroups.csv > gene_famil
ies_lost_riftia_owenia.csv
16 cut -f 1,21 gene_families_lost_riftia_owenia.csv | sed 's/Ofu
s|//g' | cut -f 2 | sed 's/, /\n/g' > gene_IDs_lost_riftia_oweni
a.txt
17
18 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
grep -w Pech | cut -f 1 > gene_families_lost_paraescarpia.txt
19 grep -f gene_families_lost_paraescarpia.txt owenia_all_GF > loss
es_paraescarpia_owenia
20 fgrep -f losses_paraescarpia_owenia ../../Orthogroups.csv > gene
_families_lost_paraescarpia_owenia.csv
21 cut -f 1,21 gene_families_lost_paraescarpia_owenia.csv | sed 's/
Ofus|//g' | cut -f 2 | sed 's/, /\n/g' > gene_IDs_lost_paraescar
pia_owenia.txt
22
23 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
grep -w Lluy | cut -f 1 > gene_families_lost_lamellibrachia.txt
```

```
24 grep -f gene_families_lost_lamellibrachia.txt owenia_all_GF > losses_lamellibrachia_owenia
25 fgrep -f losses_lamellibrachia_owenia ../../Orthogroups.csv > gene_families_lost_lamellibrachia_owenia.csv
26 cut -f 1,21 gene_families_lost_lamellibrachia_owenia.csv | sed 's/Ofus|//g' | cut -f 2 | sed 's/, /\n/g' > gene_IDs_lost_lamellibrachia_owenia.txt
27
28 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv | grep -w Siboglinidae | cut -f 1 > gene_families_lost_Siboglinidae.txt
29 grep -f gene_families_lost_Siboglinidae.txt owenia_all_GF > losses_Siboglinidae_owenia
30 fgrep -f losses_Siboglinidae_owenia ../../Orthogroups.csv > gene_families_lost_Siboglinidae_owenia.csv
31 cut -f 1,21 gene_families_lost_Siboglinidae_owenia.csv | sed 's/Ofus|//g' | cut -f 2 | sed 's/, /\n/g' > gene_IDs_lost_Siboglinidae_owenia.txt
32
33 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv | grep -w Vestimentifera | cut -f 1 > gene_families_lost_Vestimentifera.txt
34 grep -f gene_families_lost_Vestimentifera.txt owenia_all_GF > losses_Vestimentifera_owenia
35 fgrep -f losses_Vestimentifera_owenia ../../Orthogroups.csv > gene_families_lost_Vestimentifera_owenia.csv
36 cut -f 1,21 gene_families_lost_Vestimentifera_owenia.csv | sed 's/Ofus|//g' | cut -f 2 | sed 's/, /\n/g' > gene_IDs_lost_Vestimentifera_owenia.txt
37
38 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv | grep -w Vestimentifera_cl1 | cut -f 1 > gene_families_lost_Vestimentifera_cl1.txt
39 grep -f gene_families_lost_Vestimentifera_cl1.txt owenia_all_GF
```

```

> losses_Vestimentifera_cl1_owenia
40 fgrep -f losses_Vestimentifera_cl1_owenia ../../Orthogroups.csv
> gene_families_lost_Vestimentifera_cl1_owenia.csv
41 cut -f 1,21 gene_families_lost_Vestimentifera_cl1_owenia.csv | s
  ed 's/Ofus|//g' | cut -f 2 | sed 's/, /\n/g' > gene_IDs_lost_Ves
  timentifera_cl1_owenia.txt
42
43 cut -f 1,6 ../../orthofinder_ultrasensitive_stats_Jun2021.tsv |
  grep -w Vestimentifera_cl2 | cut -f 1 > gene_families_lost_Vesti
  mentifera_cl2.txt
44 grep -f gene_families_lost_Vestimentifera_cl2.txt owenia_all_GF
> losses_Vestimentifera_cl2_owenia
45 fgrep -f losses_Vestimentifera_cl2_owenia ../../Orthogroups.csv
> gene_families_lost_Vestimentifera_cl2_owenia.csv
46 cut -f 1,21 gene_families_lost_Vestimentifera_cl2_owenia.csv | s
  ed 's/Ofus|//g' | cut -f 2 | sed 's/, /\n/g' > gene_IDs_lost_Ves
  timentifera_cl2_owenia.txt

```

Various

Piecharts of composition lost gene families

```

1 03s
2 mkdir piecharts
3 cd piecharts
4 cp /data/SBCS-MartinDuranLab/03-Giacomo/data/all_together/gene_f
  amily_evolution/orthofinder_Jun2021/ultra_sensitive/Ferdi_result
  /orthofinder_ultrasensitive_stats_Jun2021.tsv ./
5 cut -f 1,6 orthofinder_ultrasensitive_stats_Jun2021.tsv | grep -
  w Ofra | cut -f 1 > gene_families_lost_osedax.txt
6 grep -f gene_families_lost_osedax.txt orthofinder_ultrasensitive
  _stats_Jun2021.tsv | cut -f 4 | sort | uniq -c > piechart_losses
  _osedax
7
8 cut -f 1,6 orthofinder_ultrasensitive_stats_Jun2021.tsv | grep -

```



```
w Dgyr | cut -f 1 > gene_families_lost_dimorphilus.txt
9 grep -f gene_families_lost_dimorphilus.txt orthofinder_ultrasens
itive_stats_Jun2021.tsv | cut -f 4 | sort | uniq -c > piechart_l
osses_dimorphilus
10
11 cut -f 1,6 orthofinder_ultrasensitive_stats_Jun2021.tsv | grep -
w Oalv | cut -f 1 > gene_families_lost_oasisia.txt
12 grep -f gene_families_lost_oasisia.txt orthofinder_ultrasensitiv
e_stats_Jun2021.tsv | cut -f 4 | sort | uniq -c > piechart_losse
s_oasisia
13
14 cut -f 1,6 orthofinder_ultrasensitive_stats_Jun2021.tsv | grep -
w Rpac | cut -f 1 > gene_families_lost_riftia.txt
15 grep -f gene_families_lost_riftia.txt orthofinder_ultrasensitive
_stats_Jun2021.tsv | cut -f 4 | sort | uniq -c > piechart_losses
_riftia
16
17 cut -f 1,6 orthofinder_ultrasensitive_stats_Jun2021.tsv | grep -
w Pech | cut -f 1 > gene_families_lost_paraescarpia.txt
18 grep -f gene_families_lost_paraescarpia.txt orthofinder_ultrasen
sitive_stats_Jun2021.tsv | cut -f 4 | sort | uniq -c > piechart_
losses_paraescarpia
19
20 cut -f 1,6 orthofinder_ultrasensitive_stats_Jun2021.tsv | grep -
w Lluy | cut -f 1 > gene_families_lost_lamellibrachia.txt
21 grep -f gene_families_lost_lamellibrachia.txt orthofinder_ultras
ensitive_stats_Jun2021.tsv | cut -f 4 | sort | uniq -c > piechar
t_losses_lamellibrachia
22
23 cut -f 1,6 orthofinder_ultrasensitive_stats_Jun2021.tsv | grep -
w Siboglinidae | cut -f 1 > gene_families_lost_siboglinidae.txt
24 grep -f gene_families_lost_siboglinidae.txt orthofinder_ultrasen
sitive_stats_Jun2021.tsv | cut -f 4 | sort | uniq -c > piechart_
losses_siboglinidae
```

```
1 ggplot(data, aes(x="", y=value, fill=group)) +  
2   geom_bar(stat="identity", width=1) +  
3   coord_polar("y", start=0)
```

```
1 cp /data/SBCS-MartinDuranLab/03-Giacomo/data/all_together/gene_f  
  amily_evolution/orthofinder_Jun2021/ultra_sensitive/Results_Jun0  
  9/Orthogroups/Orthogroups.csv ./  
2 grep "Ofra|" Orthogroups.csv | cut -f 1 > orthogroups_Ofra  
3 grep -f orthogroups_Ofra orthofinder_ultrasensitive_stats_Jun202  
  1.tsv | cut -f 4 | sort | uniq -c > piechart_Ofra  
4  
5 grep "Oalv|" Orthogroups.csv | cut -f 1 > orthogroups_Oalv  
6 grep -f orthogroups_Oalv orthofinder_ultrasensitive_stats_Jun202  
  1.tsv | cut -f 4 | sort | uniq -c > piechart_Oalv  
7  
8 grep "Rpac|" Orthogroups.csv | cut -f 1 > orthogroups_Rpac  
9 grep -f orthogroups_Rpac orthofinder_ultrasensitive_stats_Jun202  
  1.tsv | cut -f 4 | sort | uniq -c > piechart_Rpac  
10  
11 grep "Hrob|" Orthogroups.csv | cut -f 1 > orthogroups_Hrob  
12 grep -f orthogroups_Hrob orthofinder_ultrasensitive_stats_Jun202  
  1.tsv | cut -f 4 | sort | uniq -c > piechart_Hrob
```

PCA - orthogroups

input is a file generated by orthofinder: Orthogroups.GeneCount.tsv

```
1 library(stats)  
2 library(factoextra)  
3 library(R.utils)  
4
```

```
5 # Whole dataset: 47,685 orthogroups
6 orthogroups <- read.csv("/Users/giacomo/Desktop/PCA/Orthogroups.
  GeneCount.tsv", header=T, sep='\t')
7 orthogroups_selection <- data.frame()
8
9 # We subset the 21,189 orthogroups that are not species-specific
10 # or the 14,373 orthogroups that are found in at least 3 species
11 for (i in 1:nrow(orthogroups)){
12   if (sum(isZero(orthogroups[i,-c(1,30)])) < 26){
13     orthogroups_selection <- rbind(orthogroups_selection,orthogr
      oups[i,])
14   }
15 }
16
17 # We transpose the orthogroups to analyse the individuals rather
  than the
18 # orthogroups in the downstream analysis``
19
20 individuals <- t(orthogroups_selection[-c(1,30)])
21 individuals_clean <- individuals[,which(apply(individuals, 2, va
  r)!=0)]
22
23 pca_results <- prcomp(individuals_clean, scale = TRUE)
24 fviz_eig(pca_results)
25 fviz_pca_ind(pca_results, repel = TRUE) +
26   theme_classic()
```

DNA repair

list

1	A-OGG1	PTHR10242
2	A-NTHL1	PTHR43286
3	A-NEIL	PTHR22993
4	A-UNG	PTHR11264
5	A-SMUG1	PTHR13235
6	A-MUTYH	PTHR42944
7	A-MPG	PTHR10429
8	A-MBD4	PTHR15074
9	A-TDG	PTHR12159
10	A-APEX	PTHR22748
11	A-ERCC1	PTHR12749
12	A-POL	PTHR11276
13	A-HMGB1	PTHR48112
14	A-POLD3	PTHR17598
15	A-POLE3	PTHR46172
16	A-PCNA	PTHR11352
17	A-LIG	PTHR45674
18	A-FEN1	PTHR11081
19	A-PARP2	PTHR10459
20	B-RBX1	PTHR11210
21	B-CUL4B	PTHR11932
22	B-DDB1	PTHR10644
23	B-DDB2	PTHR15169
24	B-XPC	PTHR12135
25	B-HR23B	PTHR10621
26	B-CETN2	PTHR23050
27	B-CSA	PTHR46202
28	B-CSB	PTHR45629
29	B-CDK7	PTHR24056

30	B-MNAT1	PTHR12683
31	B-CCNH	PTHR10026
32	B-XPB	PTHR11274
33	B-XPD	PTHR11472
34	B-TTDA	PTHR28580
35	B-TFIIH1	PTHR12856
36	B-TFIIH2	PTHR12695
37	B-TFIIH3	PTHR12831
38	B-TFIIH4	PTHR13152
39	B-XPG	PTHR16171
40	B-XPA	PTHR10142
41	B-RPA	PTHR13989
42	B-XPF	PTHR10150
43	B-ERCC1	PTHR12749
44	B-POLD3	PTHR17598
45	B-POLE3	PTHR46172
46	B-PCNA	PTHR11352
47	B-RFC1	PTHR23389
48	B-LIGI	PTHR45674
49	C-PMS2-MLH	PTHR10073
50	C-MSH	PTHR11361
51	C-RFC1	PTHR23389
52	C-RFC	PTHR11669
53	C-PCNA	PTHR11352
54	C-EXO1	PTHR11081
55	C-RPA	PTHR13989
56	C-POLD3	PTHR17598
57	C-LIGI	PTHR45674
58	D-ATM	PTHR11139
59	D-RAD50	PTHR18867

60	D-MRE11	PTHR10139
61	D-NBS1	PTHR12162
62	D-TOPB1	PTHR13561
63	D-CTIP	PTHR15107
64	D-BARD1	PTHR24171
65	D-BRCA1	PTHR13763
66	D-BRIP1	PTHR11472
67	D-PALB2	PTHR14662
68	D-BRCA2	PTHR11289
69	D-DSS1	PTHR16771
70	D-SYCP3	PTHR19368
71	D-ABRAXAS1	PTHR31728
72	D-RAP80	PTHR15932
73	D-NPA1	PTHR15660
74	D-BRE	PTHR15189
75	D-BRCC36	PTHR10410
76	D-RPA	PTHR13989
77	D-RAD51	PTHR22942
78	D-RAD52	PTHR12132
79	D-RAD54	PTHR45629
80	D-POLD3	PTHR17598
81	D-BLM	PTHR13710
82	D-MUS81	PTHR13451
83	D-TOP3	PTHR11390
84	D-EME1	PTHR21077
85	D-RAD51B	PTHR46456
86	D-RAD51C	PTHR46239
87	D-RAD51D	PTHR46457
88	D-XRCC2	PTHR46644
89	D-XRCC3	PTHR46487

90	E-KU	PTHR12604
91	E-RAD50	PTHR18867
92	E-MRE11	PTHR10139
93	E-ARTEMIS	PTHR23240
94	E-DNAPK	PTHR11139
95	E-RAD27	PTHR11081
96	E-POL-TDT	PTHR11276
97	E-LIG	PTHR45997
98	E-XRCC4	PTHR28559
99	E-XLF	PTHR32235
100	F-MHF	PTHR22980
101	F-FANCM	PTHR14025
102	F-FAAP24	PTHR31786
103	F-TEL02	PTHR15830
104	F-ATRIP	PTHR28594
105	F-ATR	PTHR11139
106	F-WDR48	PTHR19862
107	F-USP1	PTHR24006
108	F-FANCI	PTHR21818
109	F-FANCD2	PTHR32086
110	F-FANCD2OS	PTHR31036
111	F-FANCB	PTHR28450
112	F-FAAP100	PTHR14890
113	F-FANCA	PTHR12047
114	F-FANCL	PTHR13206
115	F-FANCC	PTHR16798
116	F-FANCE	PTHR32094
117	F-FANCG	PTHR15254
118	F-FANCF	PTHR14449
119	F-UBE2T	PTHR24068

120	F-HES1	PTHR10985
121	F-MUS81	PTHR13451
122	F-EME1	PTHR21077
123	F-ERCC1	PTHR12749
124	F-XRCC4	PTHR28559
125	F-SLX1A	PTHR20208
126	F-SLX4	PTHR21541
127	F-RMI1	PTHR14790
128	F-RMI2	PTHR33962
129	F-TOP3A	PTHR11390
130	F-BLM	PTHR13710
131	F-RPA	PTHR13989
132	F-REV1	PTHR45990
133	F-REV3L	PTHR45812
134	F-POLH	PTHR45873
135	F-POLI	PTHR46404
136	F-POLK	PTHR11076
137	F-POLN	PTHR10133
138	F-BRCA1	PTHR13763
139	F-BRIP1	PTHR11472
140	F-PALB2	PTHR14662
141	F-BRCA2	PTHR11289
142	F-RAD51	PTHR22942
143	F-RAD51C	PTHR46239
144	F-FAN1	PTHR15749
145	F-PMS2-MLH	PTHR10073

```
sed -i 's/ */\t/g' DNA_repair_pantherID.txt
```

A - Base Excision Repair

B - Nucleotide Excision Repair

- C - Mismatch repair
- D - Homologous Recombination
- E - Non Homologous End - Joining
- F - Fanconi Anemia Pathway
- G - Microhomology-end joining repair pathway

code

test.sh

```
1  #!/bin/bash
2
3  species=$1
4  xls=/data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL_isoforms_annota
notations/"$species"*xls
5  output2="$species"_DNArepair_count.txt
6
7  while read line; do
8      Panther_ID=$(cut -f 2 <<< "$line")
9      echo $Panther_ID
10     count=$(fgrep "$Panther_ID" $xls | wc -l)
11     echo $count >> $output2
12 done < DNA_repair_pantherID.txt
```

```
1  cut -f1 DNA_repair_pantherID.txt > DNA_repair_IDs
2  echo "gene"'\t'"Riftia"'\t'"Oasisia"'\t'"Osedax"'\t'"Paraesc
arpia"'\t'"Lamellibrachia"'\t'"Owenia"'\t'"Capitella" > DNA_r
epair_siboglinidae.txt
3  paste DNA_repair_IDs riftia_DNArepair_count.txt oasisia_DNArepa
r_count.txt osedax_DNArepair_count.txt paraescarpia_DNArepair_co
unt.txt lamellibrachia_DNArepair_count.txt owenia_DNArepair_coun
t.txt capitella_DNArepair_count.txt >> DNA_repair_siboglinidae.t
xt
```

R plot

```
1 library(tidyverse)
2 library(dplyr)
3 library(ggplot2)
4 library(data.table)
5 library(gplots)
6 library(pheatmap)
7 library(dendextend)
8 library(factoextra)
9 library(ComplexHeatmap)
10 library(RColorBrewer)
11 library(NbClust)
12 library(scales)
13
14 #Import data in matrix format
15 A <- read.delim("~/Desktop/DNA_repair/DNA_repair_NOisoforms/A",
16 row.names=1)
17 B <- read.delim("~/Desktop/DNA_repair/DNA_repair_NOisoforms/B",
18 row.names=1)
19 C <- read.delim("~/Desktop/DNA_repair/DNA_repair_NOisoforms/C",
20 row.names=1)
21 D <- read.delim("~/Desktop/DNA_repair/DNA_repair_NOisoforms/D",
22 row.names=1)
23 E <- read.delim("~/Desktop/DNA_repair/DNA_repair_NOisoforms/E",
24 row.names=1)
25 F <- read.delim("~/Desktop/DNA_repair/DNA_repair_NOisoforms/F",
26 row.names=1)
27
28 # Option 1. 0 to 1 relative abundance/expression (54 is the high
29 est value in my dataset)
30 rescale_custom <- function(x) (x/14)
```

```
24 A_normalised <- t(apply(A, 1, rescale_custom))
25 rescale_custom <- function(x) (x/54)
26 B_normalised <- t(apply(B, 1, rescale_custom))
27 rescale_custom <- function(x) (x/11)
28 C_normalised <- t(apply(C, 1, rescale_custom))
29 rescale_custom <- function(x) (x/18)
30 D_normalised <- t(apply(D, 1, rescale_custom))
31 rescale_custom <- function(x) (x/18)
32 E_normalised <- t(apply(E, 1, rescale_custom))
33 rescale_custom <- function(x) (x/27)
34 F_normalised <- t(apply(F, 1, rescale_custom))
35
36
37 # To make 0 a different colour
38 # First create whatever gradient (e.g. RdBu)
39 heatmap_color <- colorRampPalette(brewer.pal(n = 7, name = "Red
40 s"))(1000)
41 heatmap_color[1] <- rgb(1,1,1)
42 #column_labels = c("your","labels"),
43 #row_labels = c("your","labels"))
44
45 paletteLength <- 1000
46 # to go from 0 to max.value (e.g. 1):
47 myBreaks <- c(seq(1/paletteLength, 1, length.out=floor(paletteLe
48 ngth)))
49
50 pheatmap(A_normalised,
51           cluster_rows = FALSE,
52           cluster_cols = FALSE,
53           border_color = NA,
54           color = heatmap_color,
```

```
53     height = 25,  
54     width = 20,  
55     breaks = myBreaks)  
56  
57 pheatmap(B_normalised,  
58     cluster_rows = FALSE,  
59     cluster_cols = FALSE,  
60     border_color = NA,  
61     color = heatmap_color,  
62     height = 25,  
63     width = 20,  
64     breaks = myBreaks)  
65  
66 pheatmap(C_normalised,  
67     cluster_rows = FALSE,  
68     cluster_cols = FALSE,  
69     border_color = NA,  
70     color = heatmap_color,  
71     height = 25,  
72     width = 20,  
73     breaks = myBreaks)  
74  
75 pheatmap(D_normalised,  
76     cluster_rows = FALSE,  
77     cluster_cols = FALSE,  
78     border_color = NA,  
79     color = heatmap_color,  
80     height = 25,  
81     width = 20,  
82     breaks = myBreaks)
```

```
83
84 pheatmap(E_normalised,
85           cluster_rows = FALSE,
86           cluster_cols = FALSE,
87           border_color = NA,
88           color = heatmap_color,
89           height = 25,
90           width = 20,
91           breaks = myBreaks)
92
93 pheatmap(F_normalised,
94           cluster_rows = FALSE,
95           cluster_cols = FALSE,
96           border_color = NA,
97           color = heatmap_color,
98           height = 25,
99           width = 20,
100           breaks = myBreaks)
101
```

- exported in PDF 8x6 inches

Developmental pathways

PFAM and panther tables

1st step

obtain the PFAM columns from the annotation files

```
cut -f 8 osedax_annotation_Jan2021_TrinoPantherK0.xls > PFAM_osedax
```

- for 5 sibo + owenia and capi

```
fgrep -f gene_families_originated_Vestimentifera.txt Orthogroups
```

```
_Jan2021.csv > gene_families_originated_Vestimentifera.csv
```

obtain transcript names from non redundant proteomes, this way we will get rid of the isoforms

```
1 grep ">" /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/Ctel.  
fa | sed 's/>Ctel_//g' > capitella_isoform_names  
2 fgrep -w -f capitella_isoform_names ../capitella_annotation_Feb2  
021_TrinoPantherK0.xls > capitella_isoform.xls  
3 grep ">" /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/Ofus.  
fa | sed 's/>_//g' > owenia_isoform_names  
4 fgrep -w -f owenia_isoform_names ../owenia_annotation_v250920.1_  
TrinoPantherK0.xls > owenia_isoform.xls  
5 grep ">" /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/Lluy.  
fa | sed 's/>Lluy_//g' > lamellibrachia_isoform_names  
6 fgrep -w -f lamellibrachia_isoform_names ../lamellibrachia_annot  
ation_Feb2021_TrinoPantherK0_OK.xls > lamellibrachia_isoform.xls  
7 grep ">" /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/Pech.  
fa | sed 's/>Pech_//g' | sed 's/nbis-mrna-/nbis_mrna_/g' > parae  
scarpia_isoform_names  
8 fgrep -w -f paraescarpia_isoform_names ../paraescarpia_annotatio  
n_Jun2021_TrinoPantherK0.xls > paraescarpia_isoform.xls  
9 grep ">" /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/Ofra.  
fa | sed 's/>Ofra_//g' > osedax_isoform_names  
10 fgrep -w -f osedax_isoform_names ../osedax_annotation_Jan2021_Tr  
inoPantherK0.xls > osedax_isoform.xls  
11 grep ">" /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/Oalv.  
fa | sed 's/>Oalv_//g' > oasisia_isoform_names  
12 fgrep -w -f oasisia_isoform_names ../oasisia_annotation_Jan2021_  
TrinoPantherK0.xls > oasisia_isoform.xls  
13 grep ">" /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/Rpac.  
fa | sed 's/>Rpac_//g' > riftia_isoform_names  
14 fgrep -w -f riftia_isoform_names ../riftia_annotation_Jan2021_Tr  
inoPantherK0.xls > riftia_isoform.xls
```

2nd step

First half:

PFAM_list.txt

```

1 T-box (PF00907) PF00907
2 Homeodomain (PF00046) PF0046
3 Fox (PF00250) PF00250
4 HMG-box (PF00505) PF00505
5 bHLH (PF00010) PF00010
6 bZIP (PF00170) PF00170
7 LAG (PF09271) PF09271
8 STAT (PF02864) PF02864
9 Mef2 (PF00319) PF00319
10 p53 (PF00870) PF00870
11 RHD (PF00554) PF00554
12 GATA (PF00320) PF00320
13 C2H2-Zn (PF00096) PF00096
14 GRH/LSF (PF04516) PF04516
15 Runt (PF00853) PF00853
16 zf-C4 (PF00105) PF00105
17 Myb (PF00249) PF00249
18 SMAD (PF03165) PF03165

```

Script:

test.sh

```

1 #!/bin/bash
2
3 species=$1
4 xls=/data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL_isoforms_an
  notations/"$species"*xls
5 output="$species"_FirstHalf_Development.txt
6
7 while read line; do

```

```

8   PFAM_ID=$(cut -f 2 <<< "$line")
9   echo $PFAM_ID
10  count=$(fgrep "$PFAM_ID" $xls | wc -l)
11  echo $count >> $output
12 done < PFAM_list.txt

```

Second half:

SecondHalf_list.txt

```

1  Notch (PF00066) PF00066
2  WNT (PF00110) PF00110
3  Fz1 (PF01534) PF01534
4  HH (PF01085) PF01085
5  Patched PTHR46022(prot itself) PTHR10796(domain containing)
6  TGF-b (PF00019) PF00019
7  TGF-bR PTHR23255(TGF-BETA RECEPTOR TYPE-2 and 1) PTHR14002(TRANS
  FORMING GROWTH FACTOR BETA RECEPTOR TYPE 3)
8  FGF (PF00167) PF00167
9  FGFR PTHR24416:SF131(type 1) (type 1) PTHR24416:SF130(type 2) PT
  HR24416:SF505(type 3) PTHR24416:SF343(type 4)
10 VEGF (PF00341) PF00341
11 VEGFR PTHR24416:SF49(type 2) PTHR24416:SF45(type 3) PTHR24416:SF
  390 (type 1)
12 EGF ligands TRANSFORMING GROWTH FACTOR ALPHA (PTHR10740) PTHR111
  00:SF18 PTHR11100:SF20 PTHR11100:SF7 PTHR46513:SF5
13 EGFR PTHR24416:SF91 PTHR24416:SF566

```

test_SecondHalf.sh

```

1  #!/bin/bash
2
3  species=$1
4  xls=/data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL_isoforms_an

```



```
notations/"$species"*xls
5 output="$species"_SecondHalf_Development.txt
6
7 grep "PF00066" $xls | wc -l >> $output
8 grep "PF00110" $xls | wc -l >> $output
9 grep "PF01534" $xls | wc -l >> $output
10 grep "PF01085" $xls | wc -l >> $output
11 grep "PTHR46022\|PTHR10796" $xls | wc -l >> $output
12 grep "PF00019" $xls | wc -l >> $output
13 grep "PTHR23255\|PTHR14002" $xls | wc -l >> $output
14 grep "PF00167" $xls | wc -l >> $output
15 grep "PTHR24416:SF131\|PTHR24416:SF130\|PTHR24416:SF505\|PTHR244
16:SF343" $xls | wc -l >> $output
16 grep "PF00341" $xls | wc -l >> $output
17 grep "PTHR24416:SF49\|PTHR24416:SF45\|PTHR24416:SF390" $xls | wc
-l >> $output
18 grep "PTHR10740\|PTHR11100:SF18\|PTHR11100:SF20\|PTHR11100:SF7\|
PTHR46513:SF5" $xls | wc -l >> $output
19 grep "PTHR24416:SF91\|PTHR24416:SF566" $xls | wc -l >> $output
```

3rd step

```
1 cut -f 1 PFAM_list.txt > list1
2 cut -f 1 SecondHalf_list.txt > list2
3
4 echo "gene"$'\t'"Riftia"$'\t'"Oasisia"$'\t'"Osedax"$'\t'"Lamelli
brachia"$'\t'"Owenia"$'\t'"Capitella" > Developmental_pathways_t
able_FirstHalf
5 echo "gene"$'\t'"Riftia"$'\t'"Oasisia"$'\t'"Osedax"$'\t'"Lamelli
brachia"$'\t'"Owenia"$'\t'"Capitella" > Developmental_pathways_t
able_SecondHalf
6
```

```

7 paste list1 riftia_FirstHalf_Development.txt oasisia_FirstHalf_Development.txt osedax_FirstHalf_Development.txt lamellibrachia_FirstHalf_Development.txt owenia_FirstHalf_Development.txt capitella_FirstHalf_Development.txt >> Developmental_pathways_table_FirstHalf
8 paste list2 riftia_SecondHalf_Development.txt oasisia_SecondHalf_Development.txt osedax_SecondHalf_Development.txt lamellibrachia_SecondHalf_Development.txt owenia_SecondHalf_Development.txt capitella_SecondHalf_Development.txt >> Developmental_pathways_table_SecondHalf

```

FirstHalf:

	gene	Riftia	Oasisia	Osedax	Lamellibr
	achia	Owenia	Capitella		
2	T-box (PF00907)	11	12	6	19
	8				9
3	Homeodomain (PF00046)	132	99	110	90
		151			117
4	Fox (PF00250)	28	30	25	31
	42				34
5	HMG-box (PF00505)	30	27	29	34
		23			28
6	bHLH (PF00010)	63	60	36	69
	81				61
7	bZIP (PF00170)	17	19	9	17
	16				13
8	LAG (PF09271)	1	1	1	1
	1				
9	STAT (PF02864)	2	1	1	2
	6				
10	Mef2 (PF00319)	2	2	2	2
	2				
11	p53 (PF00870)	1	1	1	1
	1				4

12	RHD (PF00554)	3	3	3	3	2
	3					
13	GATA (PF00320)	12	10	7	12	5
	9					
14	C2H2-Zn (PF00096)	243	263	95	262	
	563 142					
15	GRH/LSF (PF04516)	2	2	2	2	2
	3					
16	Runt (PF00853)	1	1	1	1	1
	1					
17	zfp-C4 (PF00105)	34	32	31	37	3
	7 37					
18	Myb (PF00249)	16	20	14	19	17
	17					
19	SMAD (PF03165)	5	5	4	5	5
	4					

SecondHalf:

1	gene	Riftia	Oasisia	Osedax	Lamellibr
	achia	Owenia	Capitella		
2	Notch (PF00066)	3	3	1	3
	5				
3	WNT (PF00110)	13	14	5	11
	12				
4	Fz1 (PF01534)	6	4	4	4
	5				
5	HH (PF01085)	2	2	1	3
	1				
6	Patched	18	16	7	38
					4
7	TGF- β	12	12	7	14
					14
8	TGF- β R	15	15	5	15
					24
9	FGF	4	4	3	4
					2
10	FGFR	1	1	2	1
					4

11	VEGF	1	1	4	1	2	2
12	VEGFR	2	1	1	1	1	2
13	EGF ligands		8	16	5	7	6
	5						
14	EGFR	3	1	1	1	2	2

```

1 # convert to log(n)
2 awk 'NR>1{for(i=1;i<=NF;i++) $i=log($i)}1' riftia_PFAM_Development.txt > riftia_PFAM_Development_logN.txt
3 awk 'NR>1{for(i=1;i<=NF;i++) $i=log($i)}1' oasisia_PFAM_Development.txt > oasisia_PFAM_Development_logN.txt
4 awk 'NR>1{for(i=1;i<=NF;i++) $i=log($i)}1' osedax_PFAM_Development.txt > osedax_PFAM_Development_logN.txt
5 awk 'NR>1{for(i=1;i<=NF;i++) $i=log($i)}1' lamellibrachia_PFAM_Development.txt > lamellibrachia_PFAM_Development_logN.txt
6 awk 'NR>1{for(i=1;i<=NF;i++) $i=log($i)}1' owenia_PFAM_Development.txt > owenia_PFAM_Development_logN.txt
7 awk 'NR>1{for(i=1;i<=NF;i++) $i=log($i)}1' capitella_PFAM_Development.txt > capitella_PFAM_Development_logN.txt
8
9 awk 'NR>1{for(i=1;i<=NF;i++) $i=log($i)}1' riftia_Panther_Development.txt > riftia_Panther_Development_logN.txt
10 awk 'NR>1{for(i=1;i<=NF;i++) $i=log($i)}1' oasisia_Panther_Development.txt > oasisia_Panther_Development_logN.txt
11 awk 'NR>1{for(i=1;i<=NF;i++) $i=log($i)}1' osedax_Panther_Development.txt > osedax_Panther_Development_logN.txt
12 awk 'NR>1{for(i=1;i<=NF;i++) $i=log($i)}1' lamellibrachia_Panther_Development.txt > lamellibrachia_Panther_Development_logN.txt
13 awk 'NR>1{for(i=1;i<=NF;i++) $i=log($i)}1' owenia_Panther_Development.txt > owenia_Panther_Development_logN.txt
14 awk 'NR>1{for(i=1;i<=NF;i++) $i=log($i)}1' capitella_Panther_Development.txt > capitella_Panther_Development_logN.txt

```

```
elopment.txt > capitella_Panther_Development_logN.txt
```

```
echo "gene"'\t'"Riftia"'\t'"Oasisia"'\t'"Osedax"'\t'"Lamelli  
brachia"'\t'"Owenia"'\t'"Capitella" > Developmental_pathways_t  
able_logN
```

```
paste list1 riftia_PFAM_Development_logN.txt oasisia_PFAM_Develo  
pment_logN.txt osedax_PFAM_Development_logN.txt lamellibrachia_P  
FAM_Development_logN.txt owenia_PFAM_Development_logN.txt capite  
lla_PFAM_Development_logN.txt >> Developmental_pathways_table_lo  
gN
```

```
paste list2 riftia_Panther_Development_logN.txt oasisia_Panther_  
Development_logN.txt osedax_Panther_Development_logN.txt lamelli  
brachia_Panther_Development_logN.txt owenia_Panther_Development_  
logN.txt capitella_Panther_Development_logN.txt >> Developmental  
_pathways_table_logN
```

result:

gene	Riftia	Oasisia	Osedax	Lamellibr
achia	Owenia	Capitella		
T-box (PF00907)	11	12	6	20
8				9
Homeodomain (PF0046)	3.09104	3.13549	3.0445	
2	2.94444	3.04452	3.04452	
Fox (PF00250)	3.3322	3.43399	3.21888	
3.46574	3.68888	3.73767		
HMG-box (PF00505)	3.58352	3.82864	3.55535	
3.3673	3.68888	3.13549		
bHLH (PF00010)	4.20469	4.29046	3.63759	
4.23411	4.29046	4.39445		
bZIP (PF00170)	2.99573	3.21888	2.30259	
2.83321	2.77259	2.77259		
LAG (PF09271)	0	0	0	1.09861
0				

9	STAT (PF02864)	0.693147	0	0	0.693147
	1.09861	1.79176			
10	Mef2 (PF00319)	0.693147	1.09861		1.09861
	0.693147	0.693147	0.693147		
11	p53 (PF00870)	1.60944	0	0	0
	1.38629	0			
12	RHD (PF00554)	1.09861	1.38629		1.38629
	1.09861	0.693147	1.09861		
13	GATA (PF00320)	2.70805	2.63906		2.07944
	2.48491	2.3979	2.30259		
14	C2H2-Zn (PF00096)	5.52943	5.63835		4.58497
	5.57215	6.44572	4.95583		
15	GRH/LSF (PF04516)	1.09861	1.38629		0.693147
	0.693147	0.693147	1.09861		
16	Runt (PF00853)	0	0	0	0
	0				
17	zf-C4 (PF00105)	3.61092	3.52636		3.49651
	3.61092	3.95124	3.61092		
18	Myb (PF00249)	2.94444	3.09104		2.77259
	2.94444	3.2581	2.89037		
19	SMAD (PF03165)	1.79176	1.79176		1.60944
	1.60944	1.94591	1.38629		
20	Notch (PF00066)	1.09861	1.09861	0	1.
	09861	1.09861	1.60944		
21	WNT (PF00110)	2.56495	2.63906		1.79176
	3.09104	2.77259	2.48491		
22	Fz1 (PF01534)	1.79176	1.38629		1.38629
	1.94591	1.38629	1.60944		
23	HH (PF01085)	0.693147	0.693147	0	1.0
	9861	0	0		
24	TGF-b (PF00019)	2.48491	2.48491		1.94591
	2.63906	2.70805	2.63906		

25	FGF (PF00167)	1.38629	1.60944	1.09861		
		1.38629	1.79176	0.693147		
26	VEGF (PF00341)	0	0	1.38629	0	
		0.693147	0.693147			
27	Patched	18	16	7	38	5
						6
28	TGF- β R	2.70805	2.70805	1.60944		2.708
		05	3.3673	2.19722		
29	FGFR	0	0	0.693147	0	0
				1.38629		
30	VEGFR	0.693147	0	0	0	0
				0.693147		
31	EGF ligands	1.79176	2.48491	1.38629		
		1.94591	1.60944	1.38629		
32	EGFR	1.09861	0	0	0	1.09861
				0.693147		

with updated result containing paraescarpia as well:

```

1 library(tidyverse)
2 library(dplyr)
3 library(ggplot2)
4 library(data.table)
5 library(gplots)
6 library(pheatmap)
7 library(dendextend)
8 library(factoextra)
9 library(ComplexHeatmap)
10 library(RColorBrewer)
11 library(NbClust)
12 library(scales)
13
14 #Import data in matrix format

```

```
15 Unique_table_log <- read.delim("~/Desktop/Developmental_pathways
    /no isoforms/Unique_table_log", row.names=1)
16
17
18 # Option 1. 0 to 1 relative abundance/expression (54 is the high
    est value in my dataset)
19
20
21
22 # To make 0 a different colour
23 # First create whatever gradient (e.g. RdBu)
24 heatmap_color <- colorRampPalette(rev(brewer.pal(n = 7, name = "
    Blues")))(100)
25 #heatmap_color[1] <- rgb(0,0,0) # here include the colour you're
    interested in
26 heatmap_color <- colorRampPalette(brewer.pal(n = 7, name = "RdB
    u"))(100)
27 heatmap_color <- heatmap_color[50:100]
28 #column_labels = c("your","labels"),
29 #row_labels = c("your","labels"))
30
31
32 # If you want common scale for different heatmaps:
33 # First define some "breaks"
34
35 pheatmap(Unique_table_log,
36           cluster_rows = FALSE,
37           cluster_cols = FALSE,
38           border_color = NA,
39           color = heatmap_color,
40           height = 25,
```



```
41         width = 20)
```

```
42
```

```
43
```

- exported in PDF 8x6 inches

Updated Plot

I will divide the plot in three parts:

- first half
- second half signals
- second half receptors

```
1 2 Notch (PF00066) Receptor
2 3 WNT (PF00110) Signal
3 4 Fz1 (PF01534) Receptor
4 5 HH (PF01085) Signal
5 6 Patched Receptor
6 7 TGF- $\beta$  Signal
7 8 TGF- $\beta$ R Receptor
8 9 FGF Signal
9 10 FGFR Receptor
10 11 VEGF Signal
11 12 VEGFR Receptor
12 13 EGF ligands Signal
13 14 EGFR Receptor
```

to divide the second half file

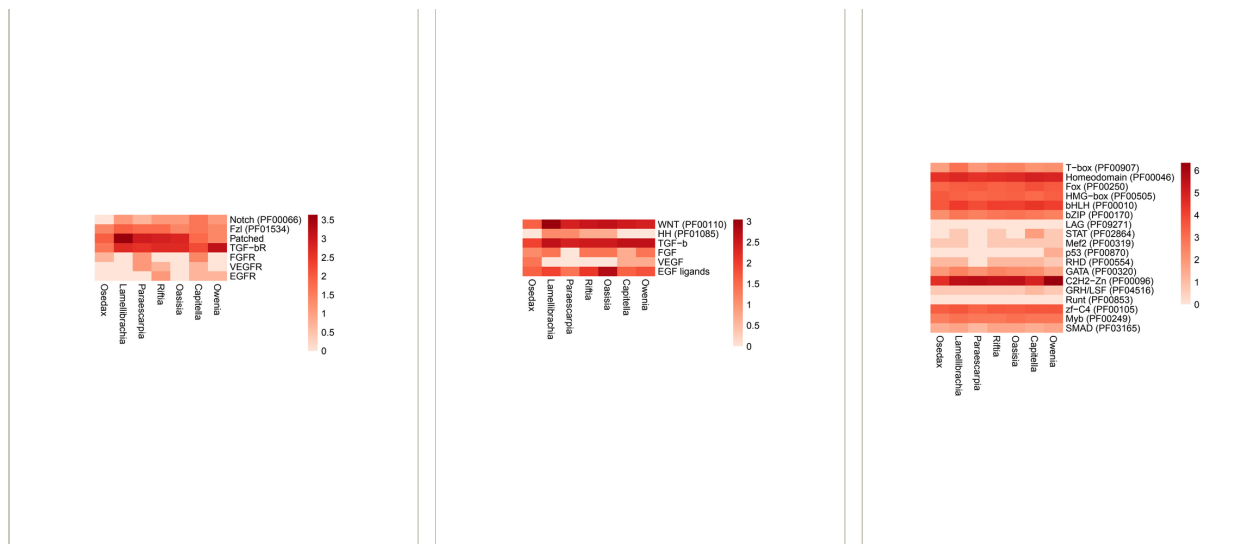
```
1 sed -e '2d;4d;6d;8d;10d;12d;14d' second_half_log > second_half_l
  og_signal
2 sed -e '3d;5d;7d;9d;11d;13d' second_half_log > second_half_log_r
  eceptor
```

```

1 library(tidyverse)
2 library(dplyr)
3 library(ggplot2)
4 library(data.table)
5 library(gplots)
6 library(pheatmap)
7 library(dendextend)
8 library(factoextra)
9 library(ComplexHeatmap)
10 library(RColorBrewer)
11 library(NbClust)
12 library(scales)
13
14 #Import data in matrix format
15 first_half_log <- read.delim("~/Desktop/Developmental_pathways/no isoforms/first_half_log_Paraescarpia_STAT=-inf", row.names=1)
16 second_half_log_signal <- read.delim("~/Desktop/Developmental_pathways/no isoforms/second_half_log_signal", row.names=1)
17 second_half_log_receptor <- read.delim("~/Desktop/Developmental_pathways/no isoforms/second_half_log_receptor", row.names=1)
18 # Option 1. 0 to 1 relative abundance/expression (54 is the highest value in my dataset)
19
20 heatmap_color <- colorRampPalette(brewer.pal(n = 7, name = "Reds"))(1000)
21
22 # If you want common scale for different heatmaps:
23 # First define some "breaks"
24
25 pheatmap(first_half_log,
```

```
26     cluster_rows = FALSE,  
27     cluster_cols = FALSE,  
28     border_color = NA,  
29     color = heatmap_color,  
30     cellheight = 10,  
31     cellwidth = 20)  
32  
33 pheatmap(second_half_log_signal,  
34     cluster_rows = FALSE,  
35     cluster_cols = FALSE,  
36     border_color = NA,  
37     color = heatmap_color,  
38     cellheight = 10,  
39     cellwidth = 20)  
40  
41 pheatmap(second_half_log_receptor,  
42     cluster_rows = FALSE,  
43     cluster_cols = FALSE,  
44     border_color = NA,  
45     color = heatmap_color,  
46     cellheight = 10,  
47     cellwidth = 20)  
48
```

- exported in PDF 8x6 inches



[PDF second_half_... • PDF document](#) [PDF second_half_... • PDF document](#) [PDF first_half_Par... • PDF document](#)

Wnt, Fz, BMP ligands & receptors

- select candidates

get transcripts id:

```
1 cd /data/scratch/btx654/developmental_pathways/WNT
2 grep "PF00110" /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL_
  isoforms_annotations/capitella* | cut -f1 | sed 's/^/Ctel_/' > W
  NT_candidates_capitella
3 module load seqtk
4 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/C
  tel.fa WNT_candidates_capitella > WNT_candidates_capitella.fa
5 grep "PF00110" /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL_
  isoforms_annotations/owenia*.xls | cut -f2 > WNT_candidates_owe
  nia
6 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/O
  fus.fa WNT_candidates_owenia > WNT_candidates_owenia.fa
7 grep "PF00110" /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL_
  isoforms_annotations/riftia* | cut -f2 | sed 's/^/Rpac_/' > WNT_
  candidates_riftia
8 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/R
  pac.fa WNT_candidates_riftia > WNT_candidates_riftia.fa
```

```

9  grep "PF00110" /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL_
   isoforms_annotations/oasisia* | cut -f2 | sed 's/^/Oalv_/' > WNT
   _candidates_oasisia
10 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/O
   alv.fa WNT_candidates_oasisia > WNT_candidates_oasisia.fa
11 grep "PF00110" /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL_
   isoforms_annotations/osedax* | cut -f2 | sed 's/^/Ofra_/' > WNT_
   candidates_osedax
12 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/O
   fra.fa WNT_candidates_osedax > WNT_candidates_osedax.fa
13 grep "PF00110" /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL_
   isoforms_annotations/lamellibrachia* | cut -f2 | sed 's/^/Lluy_
   /' > WNT_candidates_lamellibrachia
14 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/L
   luy.fa WNT_candidates_lamellibrachia > WNT_candidates_lamellibra
   chia.fa
15 grep "PF00110" /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL_
   isoforms_annotations/paraescarpia* | cut -f2 | sed 's/nbis_mrna_
   /nbis-mrna-/' | sed 's/^/Pech_/' > WNT_candidates_paraescarpia
16 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/P
   ech.fa WNT_candidates_paraescarpia > WNT_candidates_paraescarpi
   a.fa

```

- added these sequences to a fasta file containing wnt genes

```

1  "PF01534" Fzl
2  "PF01534" BMP ligands
3  "PTHR23255\|PTHR14002" BMP receptors

```

• Phylogenetic reconstruction

Step 1 - WNT phylogenetic tree

make an alignment using the txt file obtained before and the online tool [MAFFT](#)

don't select any additional options

export in fasta format and name the output:

```
wnt_sequences_MAFFT1.fasta
```

Step 2 - WNT phylogenetic tree

download and install [Jalview](#)

open the file:

```
wnt_sequences_MAFFT1.fasta
```

Cut away all the areas with no conservation. Basically we should maintain only the domain

and save the cutted file as:

```
wnt_sequences_Jalview.fasta
```

Step 3 - WNT phylogenetic tree

make an alignment again using the file “wnt_sequences_Jalview.fasta” obtained before and the online tool [MAFFT](#)

select the option “L-INS-i”

export in fasta format and name the output:

```
wnt_sequences_MAFFT2.fasta
```

Step 4 - WNT phylogenetic tree

use Jalview again. Open the file:

```
wnt_sequences_MAFFT2.fasta
```

No chopping this time, just open the file and export it to fasta as:

```
wnt_sequences_Jalview2.fasta
```

Step 5 - WNT phylogenetic tree

```
wnt_sequences_gblocks.fa
```

or use trimAl

on my personal pc

```
1 conda create -n fasttree_env
2 conda activate fasttree_env
3 conda install -c bioconda fasttree
4 conda install -c bioconda trimal
```

```
trimal -in 6a.fasta -out wnt_sequences_trimal.fa
```

Step 6 - WNT phylogenetic tree

make a tree

```
1 conda activate fasttree_env
2 FastTree wnt_sequences_trimal.fa > wnt_sequences_a.tree
```

and then open it with FigTree

PARAHOXES

```
1 GS HOMEBOX 1 (PTHR24339:SF29) GSX1
2 GS HOMEBOX 2 (PTHR47421:SF1) GSX2
3 PANCREAS/DUODENUM HOMEBOX PROTEIN 1 (PTHR45664:SF12) PDX1
4 HOMEBOX PROTEIN CDX-1 (PTHR24332:SF16)
5 CDX2 PTHR24332:SF27
6 CDX4 PTHR24332:SF15
```

parahoxes_Ofus_Ctel.fa

```
1 >CDX1_owenia
2 MVQEFGETLHYTAGTRQPTTMSLNPFLTQQGGYPQDFGPAFQISPMESQMQQWGMYAGSRS
  AAGLEEWQQAFAAQAQNAGAYAAHFNQTGHPAPMSASVNTTSPRQQNRAPFDWMKRQTYAAQPA
  AGKTRTKDKYRVVYSDHQRLLEKEFHYSRYITIRKAEALSAQALSLSERQVKIWFQNRRAKERK
  CNKKKDEQNSILGKEDSELIVSEHLTQEHVHQALSAQQHV
3 >GSX2_owenia
4 MSTSYFVDALLKKPTQMSLQRELSTAMSRQSQITHLPPPAHNHTPMLPGQPLACYPRRPSELF
  GGCCPLCIQTPGGHLICPSNAASNLSTMKHLPTSSASALSSSSGFTPRLPLAINTVSRLPSRR
  DSPSPPEYSAVDTRRIRYMNLGNIGMTRDSSDDLPSGKRIRTAFTSTQLELEREFASNMYLSR
  LRRIEIATYLNLSEKQVKIWFQNRVRKQKKEGTDEAPTHDKCRCLRTCASRNEKQDIECHGNDC
  ESVNSPSSEIDSSDINSSDISQVSSQSITKDSQIPVDING
5 >PDX1_owenia
6 MDGSNPYCSQGMVGRDSYGGHTQQSMGPYNLPACVYDTNKQTSIGLDYTSQHGMAMVDHMVEQP
```

```
MVNSIPAHSLSQPQPTPAHMQHGVNMNISNLNTNVPQQQSHPTSVPLQPPPAHQSQKPSNSNGGN
SSSSNNNNEKPLQFPWMKTTKSHARQWKAQWPGANFNIEDDNKRTRTAYTRLQLVELEKEFHYS
KYISRPRRIEIAAMLNLTERHIKIWFQNRMRKWKKDEAKRRPRPLSEEIDSKVAINTELLDKDG
AGSSPEIMTKEENPTFDDLSDSLSPSNTKPFIGTMKD
```

```
>CDX_capitella
```

```
MVQELSYLDSYSSNFPMFHHQPPSGGALRHQGPPYMHAPSATGLSGYGVQQPTNNMLMASQQYN
HSQDYSNYMDSTGRAQQTTCSPVQSSAGWSTGGMYAAGPVHTPRMDEWGNFVSQPSLQIPPNQ
HSPGTAYPPYRSPDQLQQAAGGGGGGGGGSGGFGVSPAPQACSQSPMPSQHMQPAQSPSPSSII
GAQGDGLSPGSPGAGGMRHQGQQNPRVPYDWMKKQNYQSMPLAMGKTRTKDKYRIVYSEYQKVE
LEKEYLYSKYITIQRKAELSRSIGLSERQVKIWFQNRRAKERKQKRKMEEALSTSDSSDKEHIK
VETMHEVHAPHMHSHAPH*
```

```
>GSX_capitella
```

```
MTSSTSFSVDTLLEYKGKPPKSSTSQIHQDSSPRSAIQPALFARTPLLPPPRNLFSETDGDKLLS
RMLCCPICLTSGHYGQICPLTIPVSANRLPSPYPHHKPVFSLSPLHVTSPFPRTHHMTIPTRFH
TGNGGHVMTQPQRPSERGSPPPESSPSGEGTKKGHSPLPCPDDDADESSDAVKRMRTAFSSTQL
LELEREFASNMYLSRLRRIEATYLSLSEKQVKIWFQNR RVKFKKEGA AHGSRDHPHCQCQLRS
CTSRDRKRDHVTEHTDIEVNVTDSDSDEKCL*
```

```
>PDX1_capitella
```

```
MEELDPCFPPPNHAAMFSRDFNGFLPSSNPYSTSESPSCVYDTRMGTSYNSGMVEDTYNHQYEN
PLPHQHQQHQRISRSRVYADPSVHLRSPNEHLVGMAHRAMDADNHVMHHVGMAPSHHSAPDKVK
EQGKVHFPWMKTTKSHAHQWKANWGANFQTFSENKRTRTAYTRAQLLELEKEFHFNRYITRPR
RVELAAHLNLTEQHIKIWFQNRMRKWKKDVKKRPQQSEQDGADDDVSSDVTDVKKIEPKIESV
QDEITDEITDENVNGVQSDQSL*
```

parahoxes_universal.sh

```
#!/bin/bash
#$ -wd /data/scratch/btx654/
#$ -o /data/scratch/btx654/
#$ -j y
#$ -pe smp 8
#$ -l h_vmem=10G
#$ -l h_rt=120:0:0
#$ -l highmem
```



```
10 species=$1
11 parahox_genes=parahoxes_0fus_Ctel.fa
12 parahox_genes_path=/data/scratch/btx654/H0X_genes/parahox/$parahox_genes
13
14 if [ "$species" == "riftia" ]; then
15     non_redundant_prot=Rpac.fa
16 fi
17 if [ "$species" == "oasisia" ]; then
18     non_redundant_prot=Oalv.fa
19 fi
20 if [ "$species" == "osedax" ]; then
21     non_redundant_prot=Ofra.fa
22 fi
23 if [ "$species" == "lamellibrachia" ]; then
24     non_redundant_prot=Lluy.fa
25 fi
26 if [ "$species" == "paraescarpia" ]; then
27     non_redundant_prot=Pech.fa
28 fi
29 non_redundant_prot_path=/data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/$non_redundant_prot
30
31 echo "Working on "$species
32
33 cd /data/scratch/btx654/H0X_genes/parahox
34 mkdir -p $species
35 cd $species
36 cp $non_redundant_prot_path ./
37 cp $parahox_genes_path ./
38
```

```
39 #make a diamond BLAST database of this proteome and BLAST the co
nsensi.fa.classified dataset against it, to find potential bona
fide genes. To make sure we only get the real genes, the e-value
is very stringent.
40 module load anaconda3
41 source activate diamond
42 diamond makedb --in $non_redundant_prot -d non_redundant_prot
43
44 diamond blastp -d non_redundant_prot -q $parahox_genes -o defaul
t.1e10.blastp -f 6 qseqid bitscore evalue stitle -k 25 -e 1e-10
-p 8
45 diamond blastp -d non_redundant_prot -q $parahox_genes -o ultra_
sensitive.1e10.blastp --ultra-sensitive -f 6 qseqid bitscore eva
lue stitle -k 25 -e 1e-10 -p 8
```

tblastn_osedax_assembly.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/H0X_genes/parahox/osedax
3 #$ -o /data/scratch/btx654/H0X_genes/parahox/osedax
4 #$ -j y
5 #$ -pe smp 8
6 #$ -l h_vmem=10G
7 #$ -l h_rt=120:0:0
8 #$ -l highmem
9
10 module load blast+
11 makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax
/haploidization/purge_dups/osedax_purged.fa -dbtype nucl -out os
edax_assembly
12 tblastn -db osedax_assembly -query parahoxes_0fus_Ctel.fa -out o
sedax_assembly_tblastn_out -max_target_seqs 5 -evalue 1e-10 -num
_threads 8 -outfmt 6
13 tblastn -db osedax_assembly -query parahoxes_0fus_Ctel.fa -out o
```

```
sedax_assembly_tblastn_out.html -max_target_seqs 5 -evaluate 1e-10  
-num_threads 8 -html
```

tblastn_osedax_transcriptome.sh

```
1  #!/bin/bash  
2  #$ -wd /data/scratch/btx654/H0X_genes/parahox/osedax  
3  #$ -o /data/scratch/btx654/H0X_genes/parahox/osedax  
4  #$ -j y  
5  #$ -pe smp 8  
6  #$ -l h_vmem=10G  
7  #$ -l h_rt=120:0:0  
8  #$ -l highmem  
9  
10 cp /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax/trinity/osed  
    ax_*/*.Trinity.fasta ./  
11 cat *.Trinity.fasta > combined.trinity.fasta  
12 rm *.Trinity.fasta  
13 module load blast+  
14 makeblastdb -in combined.trinity.fasta -dbtype nucl -out osedax_  
    transcriptome  
15 tblastn -db osedax_transcriptome -query parahoxes_0fus_Ctel.fa -  
    out osedax_transcriptome_tblastn_out -max_target_seqs 5 -evaluate  
    1e-10 -num_threads 8 -outfmt 6  
16 tblastn -db osedax_transcriptome -query parahoxes_0fus_Ctel.fa -  
    out osedax_transcriptome_tblastn_out.html -max_target_seqs 5 -ev  
    alue 1e-10 -num_threads 8 -html
```

HOX

```
1  module load anaconda3  
2  conda activate diamond  
3  conda install -c bioconda diamond
```

/data/home/btx654/scripts/06-other_analyses/HOX_genes
/diamond_blastp_universal.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 8
6  #$ -l h_vmem=10G
7  #$ -l h_rt=120:0:0
8  #$ -l highmem
9
10 species=$1
11 hox_genes=hox_genes_owenia.fa
12 hox_genes_path=/data/SBCS-MartinDuranLab/03-Giacomo/data/06-other_analyses/HOX_genes/$hox_genes
13
14 if [ "$species" == "riftia" ]; then
15     non_redundant_prot=Rpac.fa
16 fi
17 if [ "$species" == "oasisia" ]; then
18     non_redundant_prot=0alv.fa
19 fi
20 if [ "$species" == "osedax" ]; then
21     non_redundant_prot=0fra.fa
22 fi
23 if [ "$species" == "lamellibrachia" ]; then
24     non_redundant_prot=Lluy.fa
25 fi
26 if [ "$species" == "paraescarpia" ]; then
27     non_redundant_prot=Pech.fa
28 fi
```

```
29 non_redundant_prot_path=/data/SBCS-MartinDuranLab/03-Giacomo/NR_
    proteomes/$non_redundant_prot
30
31 echo "Working on "$species
32
33
34 mkdir -p HOX_genes
35 cd HOX_genes
36 mkdir -p $species
37 cd $species
38 cp $non_redundant_prot_path ./
39 cp $hox_genes_path ./
40
41 #make a diamond BLAST database of this proteome and BLAST the co
    nsensi.fa.classified dataset against it, to find potential bona
    fide genes. To make sure we only get the real genes, the e-value
    is very stringent.
42 module load anaconda3
43 source activate diamond
44 diamond makedb --in $non_redundant_prot -d non_redundant_prot
45
46 diamond blastp -d non_redundant_prot -q $hox_genes -o default.1e
    10.blastp -f 6 qseqid bitscore evalue stitle -k 25 -e 1e-10 -p 8
47 diamond blastp -d non_redundant_prot -q $hox_genes -o ultra_sens
    itive.1e10.blastp --ultra-sensitive -f 6 qseqid bitscore evalue
    stitle -k 25 -e 1e-10 -p 8
```

organising the output of blast

```
1 sort -u -k 1,1 ultra_sensitive.1e10.blastp > best_hit_riftia
2 sort -u -k 1,1 ultra_sensitive.1e10.blastp > best_hit_oasisia
3 sort -u -k 1,1 ultra_sensitive.1e10.blastp > best_hit_osedax
4 sort -u -k 1,1 ultra_sensitive.1e10.blastp > best_hit_lamellibra
```

```
chia
```

```
5 sort -u -k 1,1 ultra_sensitive.1e10.blastp > best_hit_paraescarpia
```

Step 0 - HOX phylogenetic tree

create a txt file with hox genes sequences from many species and the candidate genes identified by blast

```
1 cut -f 4 ultra_sensitive.1e10.blastp | sort | uniq > riftia_candidates
2 cut -f 4 ultra_sensitive.1e10.blastp | sort | uniq > oasisia_candidates
3 cut -f 4 ultra_sensitive.1e10.blastp | sort | uniq > osedax_candidates
4 cut -f 4 ultra_sensitive.1e10.blastp | sort | uniq > lamellibrachia_candidates
5 cut -f 4 ultra_sensitive.1e10.blastp | sort | uniq > paraescarpia_candidates
```

```
1 module load seqtk
2 seqtk subseq Rpac.fa riftia_candidates > riftia_candidates.fa
3 seqtk subseq Oalv.fa oasisia_candidates > oasisia_candidates.fa
4 seqtk subseq Ofra.fa osedax_candidates > osedax_candidates.fa
5 seqtk subseq Lluy.fa lamellibrachia_candidates > lamellibrachia_candidates.fa
6 seqtk subseq Pech.fa paraescarpia_candidates > paraescarpia_candidates.fa
7 cd ..
8 cat */*candidates.fa > candidates.fa
```

- then I have appended this sequences to a txt file containing the hox genes sequences of many different organisms. I used nano to edit on my personal pc a file sent by Oceane

Step 1 - HOX phylogenetic tree

make an alignment using the txt file obtained before and the online tool [MAFFT](#)

don't select any additional options

export in fasta format and name the output:

```
hoxes_sequences_MAFFT1.fasta
```

Step 2 - HOX phylogenetic tree

download and install [Jalview](#)

open the file:

```
hoxes_sequences_MAFFT1.fasta
```

Cut away all the areas with no conservation. Basically we should maintain only the domain

and save the cutted file as:

```
hoxes_sequences_Jalview.fasta
```

Step 3 - HOX phylogenetic tree

make an alignment again using the file "hoxes_sequences_Jalview.fasta" obtained before and the online tool [MAFFT](#)

select the option "L-INS-i"

export in fasta format and name the output:

```
hoxes_sequences_MAFFT2.fasta
```

Step 4 - HOX phylogenetic tree

use Jalview again. Open the file:

```
hoxes_sequences_MAFFT2.fasta
```

No chopping this time, just open the file and export it to fasta as:

```
hoxes_sequences_Jalview2.fasta
```

Step 5 - HOX phylogenetic tree

use the online tool [Gblocks](#).

Load the file "hoxes_sequences_Jalview2.fasta"

select all the 3 options for less stringent

copy paste the result in fasta format as:

```
hoxes_sequences_gblocks.fa
```

or use trimAl

on my personal pc

```
1 conda activate fasttree_env
```

```
2 conda install -c bioconda trimal
```

```
trimal -in hoxes_sequences_Jalview2.fasta -out hoxes_sequences_t  
rimal.fa
```

Step 6 - HOX phylogenetic tree

This part on Apocrita

```
scp -i ~/.ssh/id_rsa_apocrita hoxes_sequences_gblocks.fa btx654@  
login.hpc.qmul.ac.uk:/data/scratch/btx654/HOX_genes
```

raxml.sh

```
1 #!/bin/bash  
2 #$ -pe smp 5  
3 #$ -l highmem  
4 #$ -l h_vmem=10G  
5 #$ -l h_rt=240:0:0  
6 #$ -cwd  
7 #$ -j y  
8 module load raxml  
9 raxmlHPC -f a -b 476 -p 903 -x 12345 -# autoMRE -m PROTGAMMAAUTO  
-s hoxes_sequences_gblocks.fa -n hoxes_trial1.tre
```

the final output tree in newick format is:

```
hoxes_trial1.tre
```

```
1 cd /data/SBCS-MartinDuranLab/03-Giacomo/data/03-other_annotation  
s/paraescarpia/step9  
2 cp /data/SBCS-MartinDuranLab/03-Giacomo/data/03-other_annotation  
s/paraescarpia/step8/annotation_report.xls ./  
3 sort paraescarpia_Panther > panther_sorted  
4 cut -f 1 panther_sorted > IDs_panther
```



```

5 cut -f 2 annotation_report.xls | tail -n +2 > IDs_all
6 fgrep -w -v -f IDs_panther IDs_all > IDs_absentPanther ### There
  are oasisia:8632 osedax:4052 riftia:8182 genes without Panther a
  nnotation
7 awk '{print $0"\t""NO PTHR""\t""NO HIT"}' IDs_absentPanther > PA
  NTHER_nohits
8 cat panther_sorted PANTHER_nohits | sort -k 1,1 > Panther_sorted
  _allgenes
9 # now we need to remove duplicated lines from panther all genes
10 awk '!a[$1]++' Panther_sorted_allgenes > Panther_sorted_allgenes
  _noduplicates
11 ## use vim to add a header in Owenia_Panther_sorted_allgenes so
  that it matches Trinotate file
12 # #gene_id          transcript_id          sprout_Top_BLASTX_hit
  RNAMMER            prot_id            prot_coords            sprout_Top_BL
  ASTP_hit            Pfam            SignalP            TmHMM            egglog
  Kegg            gene_ontology_BLASTX            gene_ontology_BLAST
  P            gene_ontology_Pfam            transcript            peptide
13 awk 'FNR == NR { lineno[$1] = NR; next} {print lineno[$1], $0;}'
  annotation_report.xls Panther_sorted_allgenes_noduplicates | sor
  t -k 1,1n | cut -d' ' -f2- > Panther_sorted_allgenes_rightorder
14 echo "ID"$'\t'"Panther_1"$'\t'"Panther_2"$'\t'"Panther_3"$'\t'"P
  anther_4"$'\t'"Panther_5" > header
15 cat header Panther_sorted_allgenes_rightorder > Panther_sorted_a
  llgenes_rightorder_ok
16 paste annotation_report.xls Panther_sorted_allgenes_rightorder_o
  k > paraescarpia_annotation_Jun2021_TrinoPanther.xls
17 awk '{print $2}' paraescarpia_KAAS_custom_SBH.txt > only_K0numbe
  rs.txt
18 #add a line at the top of this file saying "K0_number"
19 nano only_K0numbers.txt
20 paste paraescarpia_annotation_Jun2021_TrinoPanther.xls only_K0nu
  mbers.txt > paraescarpia_annotation_Jun2021_TrinoPantherK0.xls

```

fasttree_env

on my personal computer

```
1 conda create -n fasttree_env
2 conda activate fasttree_env
3 conda install -c bioconda fasttree
```

```
1 FastTree alignment.file > tree_file
2
```

Step 7 - HOX phylogenetic tree

I have removed many siboglinids sequences from the tree that were clustering in a weird way. Now I can blast the siboglinid Hox genes against the mRNA of the 5 species in order to fill the gaps of the missing hoxes

first I need to obtain the the mRNA of the 5 species

```
scp -i ~/.ssh/id_rsa_apocrita /Volumes/5T\ hard-disk/Data/riftia
/Annotation/steps/step1/softmasking/riftia_softmasked.fa btx654@
login.hpc.qmul.ac.uk:/data/SBCS-MartinDuranLab/03-Giacomo/data/r
iftia/annotation/
```

```
1 cd /data/scratch/btx654/HOX_genes/further
2 module load anaconda3
3 source activate augustus
4
5 gffread -w riftia_mRNA.fa -g /data/SBCS-MartinDuranLab/03-Giacom
o/data/riftia/annotation/riftia_softmasked.fa /data/SBCS-MartinD
uranLab/03-Giacomo/data/riftia/annotation/New_annotation_Dec2020
/step6/riftia_annotation_v101220.gff3
6 gffread -w osedax_mRNA.fa -g /data/SBCS-MartinDuranLab/03-Giacom
o/data/osedax/annotation/softmasking/osedax_softmasked.fa /data/
SBCS-MartinDuranLab/03-Giacomo/data/osedax/annotation/New_annota
tion_Dec2020/step6/osedax_annotation_v101220.gff3
```

```
7 gffread -w oasisia_mRNA.fa -g /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/annotation/new_annotation_Nov2020/step1/softmasking/oasisia_Nov2020_softmasked.fa /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/annotation/New_annotation_Dec2020/step6/oasisia_annotation_v101220.gff3
8 gffread -w lamellibrachia_mRNA.fa -g /data/SBCS-MartinDuranLab/03-Giacomo/data/03-other_annotaions/lamellibrachia/new_multifasta_lamellibrachia.fa /data/SBCS-MartinDuranLab/03-Giacomo/data/03-other_annotaions/lamellibrachia/lamellibrachia_lociMerged_longestIsoform.gff
9 cp /data/SBCS-MartinDuranLab/03-Giacomo/data/03-other_annotaions/paraescarpia/step8/paraescarpia_mRNA.fa ./
10 gffread -y lamellibrachia_proteins.fa -g /data/SBCS-MartinDuranLab/03-Giacomo/data/03-other_annotaions/lamellibrachia/new_multifasta_lamellibrachia.fa /data/SBCS-MartinDuranLab/03-Giacomo/data/03-other_annotaions/lamellibrachia/lamellibrachia_lociMerged_longestIsoform.gff
11 cp /data/SBCS-MartinDuranLab/03-Giacomo/data/03-other_annotaions/paraescarpia/step8/paraescarpia_proteins.fa ./
```

Then I need to obtain the hox sequences to use for my search:

list

```
1 FUN_032673-T1
2 FUN_032670-T1
3 FUN_032669-T1
4 FUN_005888-T1
5 FUN_038049-T1
6 FUN_038047-T1
7 FUN_030985-T1
8 FUN_038044-T1
9 nbis-mrna-10157
10 nbis-mrna-10158
```

extract the sequences with seqtk

```
1 module load seqtk
```

```
2 seqtk subseq lamellibrachia_proteins.fa list > Hoxes_single_prot.fa
3 seqtk subseq paraescarpia_proteins.fa list >> Hoxes_single_prot.fa
4 fold -w 60 Hoxes_single_prot.fa > Hoxes_prot.fa #header should not be longer than 60 characters
5 cat /your/path/to/folder/*.fa > newname.fa
```

Make blast databases

tblastn.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/H0X_genes/further
3 #$ -o /data/scratch/btx654/H0X_genes/further
4 #$ -j y
5 #$ -pe smp 8
6 #$ -l h_vmem=10G
7 #$ -l h_rt=120:0:0
8 #$ -l highmem
9
10 species=$1
11 proteins="$species"_proteins.fa
12 database="$species"_db
13 output="$species"_tblastn_out
14
15 module load blast+
16 #makedb $mRNA -o $database
17 makeblastdb -in $mRNA -dbtype nucl -out $database
18 tblastn -db $database -query Hoxes.fa -out $output -max_target_seqs 25 -evalue 1e-10 -num_threads 8 -outfmt 6
```

to add

FUN_033852-T1 lamellibrachia

Step 8 - HOX phylogenetic tree

what I can do now is re-run single trees for all my species using the Hox genes

Oceane sent me and the candidate genes obtained from tblastn one species per time

```
1 cut -f 2 paraescarpia_tblastn_out | sort | uniq > paraescarpia_c
  andidates_list
2 sed -i 's/^/Pech_/' paraescarpia_candidates_list
3 seqtk subseq ../paraescarpia/Pech.fa paraescarpia_candidates_lis
  t > paraescarpia_candidates.fa
4 cut -f 2 riftia_tblastn_out | sort | uniq > riftia_candidates_li
  st
5 sed -i 's/^/Rpac_/' riftia_candidates_list
6 seqtk subseq ../riftia/Rpac.fa riftia_candidates_list > riftia_c
  andidates.fa
7 cut -f 2 osedax_tblastn_out | sort | uniq > osedax_candidates_li
  st
8 sed -i 's/^/Ofra_/' osedax_candidates_list
9 seqtk subseq ../osedax/Ofra.fa osedax_candidates_list > osedax_c
  andidates.fa
10 cut -f 2 oasisia_tblastn_out | sort | uniq > oasisia_candidates_
  list
11 sed -i 's/^/Oalv_/' oasisia_candidates_list
12 seqtk subseq ../oasisia/Oalv.fa oasisia_candidates_list > oasisi
  a_candidates.fa
13 cut -f 2 lamellibrachia_tblastn_out | sort | uniq > lamellibrach
  ia_candidates_list
14 sed -i 's/^/Lluy_/' lamellibrachia_candidates_list
15 seqtk subseq ../lamellibrachia/Lluy.fa lamellibrachia_candidates
  _list > lamellibrachia_candidates.fa
```

- then copy/paste *_candidates.fa at the end of the file Oceane sent with all the

Hoxes sequences

Now do step 1 to 6 for each single species

```
scp -i ~/.ssh/id_rsa_apocrita -r btx654@login.hpc.qmul.ac.uk:/data/scratch/btx654/H0X_genes/further/Oasisia_lox/*html /Users/giacomo/Desktop/H
```

```
1 cut -f 2 paraescarpia_tblastn_out | sort | uniq > paraescarpia_candidates_list
2 sed -i 's/^/Pech_/' paraescarpia_candidates_list
3 seqtk subseq ../paraescarpia/Pech.fa paraescarpia_candidates_list > paraescarpia_candidates.fa
```

Using Lox2 and 4 from Oasisia to search in lamellibrachia list

```
1 Oalv_OALVG000000019412.1
2 Oalv_OALVG000000019413.1
```

```
seqtk subseq ../../oasisia/Oalv.fa list > Hoxes_prot.fa
```

using owenia hoxes

tblastn_riftia_transcriptomes.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/H0X_genes/further/owenia_hoxes/transcriptomes_riftia
3 #$ -o /data/scratch/btx654/H0X_genes/further/owenia_hoxes/transcriptomes_riftia
4 #$ -j y
5 #$ -pe smp 8
6 #$ -l h_vmem=10G
7 #$ -l h_rt=120:0:0
8 #$ -l highmem
```

```
9
10 module load blast+
11 makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/riftia
    /trinity/riftia_crown_trinity/riftia_crown_trinity.Trinity.fasta
    -dbtype nucl -out riftia_crown
12 tblastn -db riftia_crown -query ../hox_genes_owenia.fa -out rift
    ia_crown_tblastn_out -max_target_seqs 25 -evaluate 1e-10 -num_thre
    ads 8 -outfmt 6
13 tblastn -db riftia_crown -query ../hox_genes_owenia.fa -out rift
    ia_crown_tblastn_out.html -max_target_seqs 25 -evaluate 1e-10 -num
    _threads 8 -html
14
15 makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/riftia
    /trinity/riftia_trunk_wall_trinity/riftia_trunk_wall_trinity.Tri
    nity.fasta -dbtype nucl -out riftia_trunkwall
16 tblastn -db riftia_trunkwall -query ../hox_genes_owenia.fa -out
    riftia_trunkwall_tblastn_out -max_target_seqs 25 -evaluate 1e-10 -
    num_threads 8 -outfmt 6
17 tblastn -db riftia_trunkwall -query ../hox_genes_owenia.fa -out
    riftia_trunkwall_tblastn_out.html -max_target_seqs 25 -evaluate 1e
    -10 -num_threads 8 -html
```

tblastn_oasisia_transcriptomes.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/H0X_genes/further/owenia_hoxes/trans
    criptomes_oasisia
3 #$ -o /data/scratch/btx654/H0X_genes/further/owenia_hoxes/transc
    riptomes_oasisia
4 #$ -j y
5 #$ -pe smp 8
6 #$ -l h_vmem=10G
7 #$ -l h_rt=120:0:0
8 #$ -l highmem
```

```
9
10 module load blast+
11 makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisi
a/trinity/oasisia_crown_trinity/oasisia_crown_trinity.Trinity.fa
sta -dbtype nucl -out oasisia_crown
12 tblastn -db oasisia_crown -query hox_genes_owenia.fa -out oasisi
a_crown_tblastn_out -max_target_seqs 25 -evaluate 1e-10 -num_threa
ds 8 -outfmt 6
13 tblastn -db oasisia_crown -query hox_genes_owenia.fa -out oasisi
a_crown_tblastn_out.html -max_target_seqs 25 -evaluate 1e-10 -num_
threads 8 -html
14
15 makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisi
a/trinity/oasisia_opistosoma_trinity/oasisia_opistosoma_trinity.
Trinity.fasta -dbtype nucl -out oasisia_opistosoma
16 tblastn -db oasisia_opistosoma -query hox_genes_owenia.fa -out o
asisia_opistosoma_tblastn_out -max_target_seqs 25 -evaluate 1e-10
-num_threads 8 -outfmt 6
17 tblastn -db oasisia_opistosoma -query hox_genes_owenia.fa -out o
asisia_opistosoma_tblastn_out.html -max_target_seqs 25 -evaluate 1
e-10 -num_threads 8 -html
18
19 makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisi
a/trinity/oasisia_trophosome_trinity/oasisia_trophosome_trinity.
Trinity.fasta -dbtype nucl -out oasisia_trophosome
20 tblastn -db oasisia_trophosome -query hox_genes_owenia.fa -out o
asisia_trophosome_tblastn_out -max_target_seqs 25 -evaluate 1e-10
-num_threads 8 -outfmt 6
21 tblastn -db oasisia_trophosome -query hox_genes_owenia.fa -out o
asisia_trophosome_tblastn_out.html -max_target_seqs 25 -evaluate 1
e-10 -num_threads 8 -html

1 cat *out | cut -f 2 | sort | uniq > oasisia_candidates_list
2 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/t
```



```

rinity/oasisia_crown_trinity/oasisia_crown_trinity.Trinity.fasta
oasisia_candidates_list > oasisia_crown_candidates.fa
3 sed -i 's/\s.*$//' oasisia_crown_candidates.fa
4 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/t
rinity/oasisia_opistosoma_trinity/oasisia_opistosoma_trinity.Tri
nity.fasta oasisia_candidates_list > oasisia_opistosoma_candidat
es.fa
5 sed -i 's/\s.*$//' oasisia_opistosoma_candidates.fa
6 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/t
rinity/oasisia_trophosome_trinity/oasisia_trophosome_trinity.Tri
nity.fasta oasisia_candidates_list > oasisia_trophosome_candidat
es.fa
7 sed -i 's/\s.*$//' oasisia_trophosome_candidates.fa

```

```

1 module load anaconda3
2 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
3/seqkit_env
3 seqkit concat oasisia*candidates.fa > oasisia_concat_candidates.
fa

```

tblastx_oasisia_concat.sh

```

1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/H0X_genes/further/owenia_hoxes/trans
criptomes_oasisia
3 #$ -o /data/scratch/btx654/H0X_genes/further/owenia_hoxes/transc
riptomies_oasisia
4 #$ -j y
5 #$ -pe smp 8
6 #$ -l h_vmem=10G
7 #$ -l h_rt=120:0:0
8 #$ -l highmem
9

```

```

10 module load blast+
11 #makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/oasis
ia/trinity/oasisia_crown_trinity/oasisia_crown_trinity.Trinity.f
asta -dbtype nucl -out oasisia_crown
12 tblastx -db ../../oasisia_db -query oasisia_concat_candidates.fa
-out oasisia_concat_tblastx_out -max_target_seqs 1 -evaluate 1e-10
-num_threads 8 -outfmt 6
13 tblastx -db ../../oasisia_db -query oasisia_concat_candidates.fa
-out oasisia_concat_tblastx_out.html -max_target_seqs 1 -evaluate
1e-10 -num_threads 8 -html

```

```

1 cut -f 2 oasisia_concat_tblastx_out | sort | uniq > oasisia_conc
at_candidates_IDs
2 sed -i 's/^/Oalv_/' oasisia_concat_candidates_IDs
3 seqtk subseq ../../../../oasisia/Oalv.fa oasisia_concat_candidates_
IDs > oasisia_concat_candidates_OK.fa

```

- from this last fasta I remove manually the sequences of already assigned hox genes

using paraescarpia hoxes

tblastn_oasisia_transcriptomes.sh

```

1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/H0X_genes/further/paraescarpia_hoxes
/transcriptomes_oasisia
3 #$ -o /data/scratch/btx654/H0X_genes/further/paraescarpia_hoxes/
transcriptomes_oasisia
4 #$ -j y
5 #$ -pe smp 8
6 #$ -l h_vmem=10G
7 #$ -l h_rt=120:0:0
8 #$ -l highmem
9
10 module load blast+

```

```

11 makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisi
    a/trinity/oasisia_crown_trinity/oasisia_crown_trinity.Trinity.fa
    sta -dbtype nucl -out oasisia_crown
12 tblastn -db oasisia_crown -query hox_genes_paraescarpia.fa -out
    oasisia_crown_tblastn_out -max_target_seqs 25 -evaluate 1e-10 -num
    _threads 8 -outfmt 6
13 tblastn -db oasisia_crown -query hox_genes_paraescarpia.fa -out
    oasisia_crown_tblastn_out.html -max_target_seqs 25 -evaluate 1e-10
    -num_threads 8 -html
14
15 makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisi
    a/trinity/oasisia_opistosoma_trinity/oasisia_opistosoma_trinity.
    Trinity.fasta -dbtype nucl -out oasisia_opistosoma
16 tblastn -db oasisia_opistosoma -query hox_genes_paraescarpia.fa
    -out oasisia_opistosoma_tblastn_out -max_target_seqs 25 -evaluate
    1e-10 -num_threads 8 -outfmt 6
17 tblastn -db oasisia_opistosoma -query hox_genes_paraescarpia.fa
    -out oasisia_opistosoma_tblastn_out.html -max_target_seqs 25 -ev
    alue 1e-10 -num_threads 8 -html
18
19 makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisi
    a/trinity/oasisia_trophosome_trinity/oasisia_trophosome_trinity.
    Trinity.fasta -dbtype nucl -out oasisia_trophosome
20 tblastn -db oasisia_trophosome -query hox_genes_paraescarpia.fa
    -out oasisia_trophosome_tblastn_out -max_target_seqs 25 -evaluate
    1e-10 -num_threads 8 -outfmt 6
21 tblastn -db oasisia_trophosome -query hox_genes_paraescarpia.fa
    -out oasisia_trophosome_tblastn_out.html -max_target_seqs 25 -ev
    alue 1e-10 -num_threads 8 -html

```

```

1 cat *out | cut -f 2 | sort | uniq > oasisia_candidates_list
2 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/t
    rinity/oasisia_crown_trinity/oasisia_crown_trinity.Trinity.fasta
    oasisia_candidates_list > oasisia_crown_candidates.fa

```

```
3 sed -i 's/\s.*$//' oasisia_crown_candidates.fa
4 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/trinity/oasisia_opistosoma_trinity/oasisia_opistosoma_trinity.Trinity.fasta oasisia_candidates_list > oasisia_opistosoma_candidates.fa
5 sed -i 's/\s.*$//' oasisia_opistosoma_candidates.fa
6 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/trinity/oasisia_trophosome_trinity/oasisia_trophosome_trinity.Trinity.fasta oasisia_candidates_list > oasisia_trophosome_candidates.fa
7 sed -i 's/\s.*$//' oasisia_trophosome_candidates.fa
```

```
1 module load anaconda3
2 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/seqkit_env
3 seqkit concat oasisia*candidates.fa > oasisia_concat_candidates.fa
```

tblastx_oasisia_concat.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/H0X_genes/further/paraescarpia_hoxes/transcriptomes_oasisia
3 #$ -o /data/scratch/btx654/H0X_genes/further/paraescarpia_hoxes/transcriptomes_oasisia
4 #$ -j y
5 #$ -pe smp 8
6 #$ -l h_vmem=10G
7 #$ -l h_rt=120:0:0
8 #$ -l highmem
9
10 module load blast+
11 #makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/oasis
```

```

ia/trinity/oasisia_crown_trinity/oasisia_crown_trinity.Trinity.f
asta -dbtype nucl -out oasisia_crown
12 tblastx -db ../../oasisia_db -query oasisia_concat_candidates.fa
   -out oasisia_concat_tblastx_out -max_target_seqs 1 -evalue 1e-10
   -num_threads 8 -outfmt 6
13 tblastx -db ../../oasisia_db -query oasisia_concat_candidates.fa
   -out oasisia_concat_tblastx_out.html -max_target_seqs 1 -evalue
   1e-10 -num_threads 8 -html

```

```

1 cut -f 2 oasisia_concat_tblastx_out | sort | uniq > oasisia_conc
  at_candidates_IDs
2 sed -i 's/^/Oalv_/' oasisia_concat_candidates_IDs
3 seqtk subseq ../../oasisia/Oalv.fa oasisia_concat_candidates_
  IDs > oasisia_concat_candidates_OK.fa

```

- from this last fasta I remove manually the sequences of already assigned hox genes

```

1 fgrep -w -f list_test lamellibrachia_annotation_Feb2021_TrinoPan
  therKO_OK.xls > IDs_absentPanther
2 cut -f 2,20 lamellibrachia_annotation_Feb2021_TrinoPantherKO_OK.
  xls | fgrep

```

single Hoxes

blast against final assemblies using single hoxes from different species

hox1.fa

```

1 >Owenia_fusiformis_Hox1
2 MNSASDYTICNLDNNTYSNNFTTDTAPYSCYANINNAGIESDSYRGGYSENNLSHHHHHHHHQH
  QHPQQLSVEAHLSGTPHNTSLSLNVYQHGTHSYPIQPSPQGNGFYEDAMIPGGELAECNAYPYP
  DNSPTSHQIAYQASEHHHIQQPPQQCDTGQQQQQSPVAQYKWMQVKRNVKPKVTDYKLTDFTY
  VNPGNNLGRNTFTNKQLTELEKEFHFNKYLTRARRIEIAAALGLNEVQVKIWFQNRMRMKQKKRL
  KENKLAQTVNGEHENGDDLSPGPTTPTSTDEQIS
3 >Capitella_tellata_Lab
4 MAGEYTLCLNDNHTYTSPYNGTEGANYNNGYTGAEYGVHHGAGPGPPGASLELHSPAGLGYGEAG

```

```

GMCDGEPAHSQAVFLHSDGQGYGAIACGGGSSAPQQHQGYAPHAGYYGHGMTFNGGIADMPPH
GLHSNGGYLGYPDPTNCNPLLGNPNASNTGYCLSPHEHVG LSSSPGSDQGPVTTYK WMTVKRG T
PKTSKAPGAGDFSVFAGQPNMGR TNFTNKQLTELEKEFHFNKYLTRARRIEIAASLGLNETQVK
IWFQNR RMKQKKRLKENTSTTPVSDSSQDGISGDLNEEAS

```

```
5 >Rpac_RPACG00000019294.1
```

```
6 MTALFGRLRQEQSIHCDLLAAGFRRLPTTAKHSEYTSYSHGAVFGSGGGAGNGTSPGGAVSAAG
QPNLGRTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLGLNETQVKIWFQNR RMKQKKRLKEG
HVCGSTNRMNDDAKSDIASLQQTADAIS
```

```
7 >Oalv_OALVG00000019490.1
```

```
8 MNFYAARSRGQWSSPVSNIYSIVSAKHTEYTSYNHGT VYGSAPGAGNGATPGGAVSAAGHPNLG
RTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLGLNETQVKIWFQNR RMKQKKRLKEGHVCGS
TNRMNDDAKSDIASLQQTADAIS
```

```
9 >Pech_nbis-mrna-10151
```

```
10 KHSEYTSYSHGAPYGS GGGVGN GATPGGAVSAVGQPNLGRTNFTNKQLTELEKEFHFNRYLTRA
RRIEIAASLGLNETQVKIWFQNR RMKQKKRLKEGHVCGLTNRMNDDAKSDMASLQQTADAIS
```

```
11 >Lluy_FUN_032673-T1
```

```
12 MNKCPLL SHHSSSPFDGTAGTYDSVTCPV TAGTKRPVRPRRSGREMSNRRFSFASKHSEYTSYN
HAGVYGN GGGGNGATPGGAVSAAGQPNLGRTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLG
LNETQVKIWFQNR RMKQKKRLKEGHVCGSTSRMND DAKSEISSLQHTADAIS
```

tblastn_osedax_assembly.sh

```

1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/H0X_genes/further/assembly/Hox1
3  #$ -o /data/scratch/btx654/H0X_genes/further/assembly/Hox1
4  #$ -j y
5  #$ -pe smp 8
6  #$ -l h_vmem=10G
7  #$ -l h_rt=120:0:0
8  #$ -l highmem
9
10 module load blast+
11 makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax
    /haploidization/purge_dups/osedax_purged.fa -dbtype nucl -out os

```

```
osedax_assembly
```

```
12 tblastn -db osedax_assembly -query hox1.fa -out osedax_assembly_
tblastn_out -max_target_seqs 5 -evalue 1e-10 -num_threads 8 -out
fmt 6
13 tblastn -db osedax_assembly -query hox1.fa -out osedax_assembly_
tblastn_out.html -max_target_seqs 5 -evalue 1e-10 -num_threads 8
-html
```

```
14
```

```
1 >osedax_1_owe
2 PQPTDFHPGEGFGFEQKRTRQTYTRYQTLELEKEFHYNRYLTRRRRIEIAHSLGLSERQIKIWFQ
NRRMKWKKENNLPKLTGPNGNDQPADSTPV
3 >osedax_2_owe
4 RTSYTRHQTLELEKEFHFNRYLTRRRRIEIAHMLTLTERQIKIWFQNRRMKWKK
5 >osedax_3_owe
6 RTAYTSAQLVELEKEFHFNRYLCRPRRIEMASLLSLSERQIKIWFQNRRMKFKK
7 >osedax_4_owe
8 RTAYTRHQVLELEKEFHFNRYLTRRRRIEIAHTLCLTERQIKIWFQNRRMKWKK
9 >osedax_5_owe
10 GRQTYSTRYQTLELEKEFQFNHYLTRKRRIEIAHVLCCLTERQIKIWFQNRRMKLKKEKQQIKDLN
DITRREHDLSPPL
11 >osedax_1_cap
12 RNPQPTDFHPGEGFGFEQKRTRQTYTRYQTLELEKEFHYNRYLTRRRRIEIAHSLGLSERQIKIW
FQNRRMKWKK
13 >osedax_2_cap
14 RTSYTRHQTLELEKEFHFNRYLTRRRRIEIAHMLTLTERQIKIWFQNRRMKWKK
15 >osedax_3_cap
16 ARTAYTSAQLVELEKEFHFNRYLCRPRRIEMASLLSLSERQIKIWFQNRRMKFKK
17 >osedax_4_cap
18 RTAYTRHQVLELEKEFHFNRYLTRRRRIEIAHTLCLTERQIKIWFQNRRMKWKK
19 >osedax_5_cap
```

```
20 SSQRRRGRQTYSTRYQTLELEKEFQFNHYLTRKRRIEIAHVLCLTERQIKIWFQNRRMKLKK
21 >osedax_1_rif
22 RQTYTRYQTLELEKEFHYNRYLTRRRRIEIAHSLGLSERQIKIWFQNRRMKWKKENNLPKLTGP
   NGNDQPAD
23 >osedax_2_rif
24 RTSYTRHQTLELEKEFHFNRYLTRRRRIEIAHMLTLTERQIKIWFQNRRMKWKK
25 >osedax_3_rif
26 RTAYTSAQLVELEKEFHFNRYLCRPRRIEMASLLSLSERQIKIWFQNRRMKFKKEQRGGIGVVG
   S
27 >osedax_4_rif
28 SNNPRRLRTAYTNTQLLELEKEFHFNKYLCRPRRIEIASTLDLTERQV
29 >osedax_5_rif
30 SRTAYTRHQVLELEKEFHFNRYLTRRRRIEIAHTLCLTERQIKIWFQNRRMKWKK
31 >osedax_1_oas
32 RTAYTSAQLVELEKEFHFNRYLCRPRRIEMASLLSLSERQIKIWFQNRRMKFKKEQRGGIGVVG
   S
33 >osedax_2_oas
34 SNNPRRLRTAYTNTQLLELEKEFHFNKYLCRPRRIEIASTLDLTERQV
35 >osedax_3_oas
36 RQTYTRYQTLELEKEFHYNRYLTRRRRIEIAHSLGLSERQIKIWFQNRRMKWKKENNLPKLTGP
   NGNDQPAD
37 >osedax_4_oas
38 RTSYTRHQTLELEKEFHFNRYLTRRRRIEIAHMLTLTERQIKIWFQNRRMKWKK
39 >osedax_5_oas
40 RTAYTRHQVLELEKEFHFNRYLTRRRRIEIAHTLCLTERQIKIWFQNRRMKWKK
41 >osedax_1_par
42 RQTYTRYQTLELEKEFHYNRYLTRRRRIEIAHSLGLSERQIKIWFQNRRMKWKKENNLPKLTGP
   NGNDQPAD
43 >osedax_2_par
44 RTSYTRHQTLELEKEFHFNRYLTRRRRIEIAHMLTLTERQIKIWFQNRRMKWKK
45 >osedax_3_par
```



```
46 RTAYTSAQLVELEKEFHFNRYLCRPRRIEMASLLSLSERQIKIWFQNRRMKFKKEQRGG
47 >osedax_4_par
48 PRRLRTAYTNTQLLELEKEFHFNKYLCRPRRIEIASTLDLTERQV
49 >osedax_5_par
50 RTAYTRHQVLELEKEFHFNRYLTRRRRIEIAHTLCLTERQIKIWFQNRRMKWKK
51 >osedax_1_lam
52 RQTYTRYQTLELEKEFHYNRYLTRRRRIEIAHSLGLSERQIKIWFQNRRMKWKK
53 >osedax_2_lam
54 RTSYTRHQVLELEKEFHFNRYLTRRRRIEIAHMLTLTERQIKIWFQNRRMKWKK
55 >osedax_3_lam
56 RTAYTSAQLVELEKEFHFNRYLCRPRRIEMASLLSLSERQIKIWFQNRRMKFKKEQRGGIGVVG
57 S
58 >osedax_4_lam
59 SNNPRRLRTAYTNTQLLELEKEFHFNKYLCRPRRIEIASTLDLTERQV
60 >osedax_5_lam
61 RTAYTRHQVLELEKEFHFNRYLTRRRRIEIAHTLCLTERQIKIWFQNRRMKWKK
```

tblastn_osedax_transcriptome.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/H0X_genes/further/assembly/Hox1
3 #$ -o /data/scratch/btx654/H0X_genes/further/assembly/Hox1
4 #$ -j y
5 #$ -pe smp 8
6 #$ -l h_vmem=10G
7 #$ -l h_rt=120:0:0
8 #$ -l highmem
9
10 cp /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax/trinity/osed
11 ax_*/*.Trinity.fasta ./
12 cat *.Trinity.fasta > combined.trinity.fasta
13 rm *.Trinity.fasta
```

```
13 module load blast+
14 makeblastdb -in combined.trinity.fasta -dbtype nucl -out osedax_
transcriptome
15 tblastn -db osedax_transcriptome -query hox1.fa -out osedax_tran
scriptome_tblastn_out -max_target_seqs 5 -evalue 1e-10 -num_thre
ads 8 -outfmt 6
16 tblastn -db osedax_transcriptome -query hox1.fa -out osedax_tran
scriptome_tblastn_out.html -max_target_seqs 5 -evalue 1e-10 -num
_threads 8 -html
```

```
1 >osedax_1_owe
2 VATYKWMTVKRNAPKTVKQTPQSSDYNGNSSTTSGACCASGSHFRSSPLSPSHSPSSIGSGCGG
GNLSGNGLPPNLGRTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMK
QKKRLKEGHAAQWTVDE
3 >osedax_2_owe
4 NLGRTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGHA
AQ
5 >osedax_3_owe
6 NLGRTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGHA
AQWTVDE
7 >osedax_4_owe
8 GSANFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKE
9 >osedax_5_owe
10 RARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKE
11 >osedax_1_cap
12 VATYKWMTVKRNAPKTVKQTPQSSDYNGNSSTTSGACCASGSHFRSSPLSPSHSPSSIGSGCGG
GNLSGNGLPPNLGRTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMK
QKKRLKE
13 >osedax_2_cap
14 GNLSGNGLPPNLGRTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMK
QKKRLKE
15 >osedax_3_cap
16 PNLGRTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKE
```

```
17 >osedax_4_cap
18 GSANFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKE
19 >osedax_5_cap
20 RARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKE
21 >osedax_1_rif
22 PNLGRTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGH
  AAQWT
23 >osedax_2_rif
24 PNLGRTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGH
  A
25 >osedax_3_rif
26 PNLGRTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRL
27 >osedax_4_rif
28 GSANFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGH
29 >osedax_5_rif
30 RARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGH
31 >osedax_1_oas
32 PNLGRTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGH
  A
33 >osedax_2_oas
34 PNLGRTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGH
  A
35 >osedax_3_oas
36 PNLGRTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGH
  A
37 >osedax_4_oas
38 GSANFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGH
39 >osedax_5_oas
40 RARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGH
41 >osedax_1_par
42 PNLGRTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGH
```

```
AAQWT
43 >osedax_2_par
44 PNLGRNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGH
AAQWT
45 >osedax_3_par
46 PNLGRNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGH
AAQWT
47 >osedax_4_par
48 GSNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGHDVL
K*RRPCALTRR
49 >osedax_5_par
50 RTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGHAAQW
T
51 >osedax_1_lam
52 PNLGRNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGH
AAQWT
53 >osedax_2_lam
54 RTNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGHAAQW
T
55 >osedax_3_lam
56 GSNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRRMKQKKRLKEGHAAQ
WT
57 >osedax_4_lam
58 RTAYTSAQLVELEKEFHFNRYLCRPRRIEMASLLSLSERQIKIWFQNRRMKFKKEQRGBGIGVVG
S
59 >osedax_5_lam
60 RTSYTRHQTLELEKEFHFNRYLTRRRRIEIAHMLTLTERQIKIWFQNRRMKWKK
```

lox2.fa

```
1 >Owenia_fusiformis_Lox2
2 MSSYFPQGQAGDMGPHDAGSSAVSEGSFNRESCTSTDFKSPGYVAPGSNYNDFTCRMPAAFQSR
TGEFRNGLPNNNFLASQYGQSGLGQGFADIDCGLGTATSLSHCSPVSPPPRVTPFYPWMSIVG
PNSNQRRRGRQTYTRFQTLELEKEFKFNRYLTRRRRIELSHMLCLTERQIKIWFQNRR
```

```

3 >Capitella_tellata_Lox2
4 MSYFNSESSTRLNGPSSVEDAGGCSLTPTVDAGLSTPPSSRLSEPGQTPTPADSVRIVSSQGYV
  PSQTHFQDYHCGTGSGISRMYESYNQHVHSNNNYLYNASQGGHLAAAAAAAAAAGGSGQPYMDLS
  VPLNCMPGSGGIPGCGGRMPGAGMNGPMYPWMSIVGPNSNQRRRGRQTYTRYQTLELEKEFKFN
  RYLTRRRRIELSHMLCLTERQIKIWFQNRRMKEKKEIQAIKELNEKEKTKPNSVPNPPTTVVD
5
6 >Oalv_OALVG00000019413.1
7 LPLALPRSYPSGITATAVAVAVAATATRREMSSFFSDHDRVEHQLRLARHAVGATAGDGPGLLG
  GADHMTSVRPPDDCGSNCAMNCRPTSGLCAPVSTFQDFACSIPIVFQSRQASADYNGGYLYSQPP
  PSSSASLSSSMTGAITAPTAHESQCHPGGQHSQTLLESHSGLAGINCGALGAANINNCGARMPP
  PQSLTAPVMYPWMSIVGPNSNQRRRGRQTYTRYQTLELEKEFKFNRYLTRRRRIELSHMLCLTE
  RQIKIWFQNRRMKEKKEIQAIKELNDKEKAKTTSATVMPSAK
8
9 >Pech_nbis-mrna-10158
10 PNSNQRRRGRQTYTRYQTLELEKEFKFNRYLTRRRRIELSHMLCLTERQIKIWFQNRRMKEKKE
  IQAIKELNEKEKMKGTSTTVLPYAK

```

tblastn_osedax_assembly.sh

```

1 #!/bin/bash
2
3 # $ -wd /data/scratch/btx654/H0X_genes/further/assembly/Lox2
4 # $ -o /data/scratch/btx654/H0X_genes/further/assembly/Lox2
5 # $ -j y
6 # $ -pe smp 8
7 # $ -l h_vmem=10G
8 # $ -l h_rt=120:0:0
9 # $ -l highmem
10
11 module load blast+
12 makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax
  /haploidization/purge_dups/osedax_purged.fa -dbtype nucl -out os
  edax_assembly
13 tblastn -db osedax_assembly -query lox2.fa -out osedax_assembly_
  tblastn_out -max_target_seqs 5 -evalue 1e-10 -num_threads 8 -out
  fmt 6
14
15 tblastn -db osedax_assembly -query lox2.fa -out osedax_assembly_

```

```
tblastn_out.html -max_target_seqs 5 -evalue 1e-10 -num_threads 8  
-html
```

```
1 >osedax_1_owe  
2 YSLLLNVGPNSSQRRRGRQTYSRYQTLELEKEFQFNHYLTRKRRIEIAHVLCLTERQIKIWFQ  
  RR  
3 >osedax_2_owe  
4 PQPTDFHPGEFGFEQKRTRQTYTRYQTLELEKEFHYNRYLTRRRRIEIAHSLGLSERQIKIWFQ  
  NRR  
5 >osedax_3_owe  
6 KRTRTSYTRHQTLELEKEFHFNRYLTRRRRIEIAHMLTLTERQIKIWFQNRR  
7 >osedax_4_owe  
8 GFNGVDSKRSRTAYTRHQVLELEKEFHFNRYLTRRRRIEIAHTLCLTERQIKIWFQNRR  
9 >osedax_5_owe  
10 KRARTAYTSAQLVELEKEFHFNRYLCRPRRIEMASLLSLSERQIKIWFQNRR  
11 >osedax_6_owe  
12 LSAESARQRKKRKPYTRYQTIMLEEEFKRNSYITRQKRWEISCKLQLSERQVKVWFQNRR  
13 >osedax_1_cap  
14 YSLLLNVGPNSSQRRRGRQTYSRYQTLELEKEFQFNHYLTRKRRIEIAHVLCLTERQIKIWFQ  
  RRMKLKKEKQQIKDLNDITRREHDLSP  
15 >osedax_1_cap  
16 LELEKEFHFNRYLTRRRRIEIAHMLTLTERQIKIWFQNRRMKWKKEHKA  
17 >osedax_2_cap  
18 LELEKEFHYNRYLTRRRRIEIAHSLGLSERQIKIWFQNRRMKWKKE  
19 >osedax_3_cap  
20 LELEKEFHFNRYLTRRRRIEIAHTLCLTERQIKIWFQNRRMKWKKD  
21 >osedax_4_cap  
22 VELEKEFHFNRYLCRPRRIEMASLLSLSERQIKIWFQNRRMKFKKE  
23 >osedax_1_oas  
24 YSLLLNVGPNSSQRRRGRQTYSRYQTLELEKEFQFNHYLTRKRRIEIAHVLCLTERQIKIWFQ  
  RRMKLKKEKQQIKDLND
```

```
25 >osedax_2_oas
26 LELEKEFHFNRYLTRRRRIEIAHMLTLTERQIKIWFQNRRMKWKKEHKA
27 >osedax_3_oas
28 LELEKEFHYNRYLTRRRRIEIAHSLGLSERQIKIWFQNRRMKWKKE
29 >osedax_4_oas
30 LELEKEFHFNRYLTRRRRIEIAHTLCLTERQIKIWFQNRRMKWKKD
31 >osedax_5_oas
32 VELEKEFHFNRYLCRPRRIEMASLLSLSERQIKIWFQNRRMKFKKE
33 >osedax_1_par
34 PNSSQRRRGRQTYSTRYQTLLELEKEFQFNHYLTRKRRIEIAHVLCLTERQIKIWFQNRRMKLKKE
  KQQIKDLND
35 >osedax_2_par
36 LELEKEFHFNRYLTRRRRIEIAHTLCLTERQIKIWFQNRRMKWKKD
37 >osedax_3_par
38 LELEKEFHFNRYLTRRRRIEIAHMLTLTERQIKIWFQNRRMKWKKEHKA
39 >osedax_4_par
40 LELEKEFHYNRYLTRRRRIEIAHSLGLSERQIKIWFQNRRMKWKKE
41 >osedax_5_par
42 VELEKEFHFNRYLCRPRRIEMASLLSLSERQIKIWFQNRRMKFKKE
```

tblastn_osedax_transcriptome.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/H0X_genes/further/assembly/Lox2
3 #$ -o /data/scratch/btx654/H0X_genes/further/assembly/Lox2
4 #$ -j y
5 #$ -pe smp 8
6 #$ -l h_vmem=10G
7 #$ -l h_rt=120:0:0
8 #$ -l highmem
9
10 cp /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax/trinity/osed
```

```
ax_*/*.Trinity.fasta ./
11 cat *.Trinity.fasta > combined.trinity.fasta
12 rm *.Trinity.fasta
13 module load blast+
14 makeblastdb -in combined.trinity.fasta -dbtype nucl -out osedax_
transcriptome
15 tblastn -db osedax_transcriptome -query lox2.fa -out osedax_tran
scriptome_tblastn_out -max_target_seqs 5 -evalue 1e-10 -num_thre
ads 8 -outfmt 6
16 tblastn -db osedax_transcriptome -query lox2.fa -out osedax_tran
scriptome_tblastn_out.html -max_target_seqs 5 -evalue 1e-10 -num
_threads 8 -html
```

```
1 >osedax_1_owe
2 ASEPSPNVMPWMSIVGPNSNQRRRGRQTYTRYQTLELEKEFKYNRYLTRRRRIELSHTLCLT
ERQIKIWFQNRR
3 >osedax_2_owe
4 GRQTYTRYQTLELEKEFKYNRYLTRRRRIELSHTLCLTERQIKIWFQNRR
5 >osedax_1_cap
6 YPWMSIVGPNSNQRRRGRQTYTRYQTLELEKEFKYNRYLTRRRRIELSHTLCLTERQIKIWFQN
RRMKEKKEIQAIKELNAKEQ
7 >osedax_2_cap
8 LELEKEFKYNRYLTRRRRIELSHTLCLTERQIKIWFQNRRMKEKKEIQAIKELNAKEQ
9 >osedax_1_oas
10 NSDNGMASEPPSPNVMPWMSIVGPNSNQRRRGRQTYTRYQTLELEKEFKYNRYLTRRRRIELS
HTLCLTERQIKIWFQNRRMKEKKEIQAIKELNAKEQQST
11 >osedax_2_oas
12 WYLV*VGPNSNQRRRGRQTYTRYQTLELEKEFKYNRYLTRRRRIELSHTLCLTERQIKIWFQNRR
MKEKKEIQAIKELNAKEQQST
13 >osedax_1_par
14 PNSNQRRRGRQTYTRYQTLELEKEFKYNRYLTRRRRIELSHTLCLTERQIKIWFQNRRMKEKKE
IQAIKELNAKEQ
```



```
15 >osedax_2_par
16 LELEKEFKYNRYLTRRRRIELSHTLCLTERQIKIWFQNRRMKEKKEIQAIKELNAKEQ
```

antp.fa

```
1 >Owenia_fusiformis_Antp
2 MSYYHNGSYMTTEHFAGHNSPMTTNYQNSPRTATIYDDTSQPAYPRFPPYDRLDIRPIQSNQQP
3 QGGYYNQNTIARDNRDDYSHTNGQQLSPIQRYSSCKISDDTVESYLAGAVADHTTPLYENNNSP
4 PLQTSISPPQPQAPQTETPNQNQSQNQNSTQQQIPIYPWMRSQFGPDRKRGRQTYTRFQTLELE
5 KEFHFNKYLTRRRRIEIAHSLCLTERQIKIWFQNRRMKWKKENKQIEALKSPESDEKSEPPSPS
6 SPTEDDDELKKEKEDDDGDKPTTPDTL
7
8 >Capitella_tellata_Antp
9 MNNWTKTKMLRIQTKLSKKNQLKAGPERKRGRTYTRYQTLELEKEFHFNRYLTRRRRIEIAHA
10 LCLTERQIKIWFQNRRMKWKKENRQIEVLRQHTDDDLDFR
```

tblastn_osedax_both.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/H0X_genes/further/assembly/Antp/osed
3 ax
4 #$ -o /data/scratch/btx654/H0X_genes/further/assembly/Antp/oseda
5 x
6 #$ -j y
7 #$ -pe smp 8
8 #$ -l h_vmem=10G
9 #$ -l h_rt=120:0:0
10 #$ -l highmem
11
12 module load blast+
13 makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax
14 /haploidization/purge_dups/osedax_purged.fa -dbtype nucl -out os
15 edax_assembly
16
17 tblastn -db osedax_assembly -query antp.fa -out osedax_assembly_
18 tblastn_out -max_target_seqs 5 -evalue 1e-10 -num_threads 8 -out
19 fmt 6
```

```
13 tblastn -db osedax_assembly -query antp.fa -out osedax_assembly_
    tblastn_out.html -max_target_seqs 5 -evalue 1e-10 -num_threads 8
    -html
14
15 cp /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax/trinity/osed
    ax_*/*.Trinity.fasta ./
16 cat *.Trinity.fasta > combined.trinity.fasta
17 rm *.Trinity.fasta
18 makeblastdb -in combined.trinity.fasta -dbtype nucl -out osedax_
    transcriptome
19 tblastn -db osedax_transcriptome -query antp.fa -out osedax_tran
    scriptome_tblastn_out -max_target_seqs 5 -evalue 1e-10 -num_thre
    ads 8 -outfmt 6
20 tblastn -db osedax_transcriptome -query antp.fa -out osedax_tran
    scriptome_tblastn_out.html -max_target_seqs 5 -evalue 1e-10 -num
    _threads 8 -html
```

```
1 >osedax_1_assembly_owe
2 GEFQFEQKRTRQTYTRYQTLELEKEFHYNRYLTRRRRIEIAHSLGLSERQIKIWFQNRRMKWKK
    EN
3 >osedax_2_assembly_owe
4 GIDSKRTRTSYTRHQTLLELEKEFHFNRYLTRRRRIEIAHMLTLTERQIKIWFQNRRMKWKKKEHK
5 >osedax_3_assembly_owe
6 GVDSKRSRTAYTRHQVLELEKEFHFNRYLTRRRRIEIAHTLCLTERQIKIWFQNRRMKWKKDHK
7 >osedax_1_assembly_cap
8 GFEQKRTRQTYTRYQTLELEKEFHYNRYLTRRRRIEIAHSLGLSERQIKIWFQNRRMKWKKENN
    LPKLTGPNQND
9 >osedax_2_assembly_cap
10 GIDSKRTRTSYTRHQTLLELEKEFHFNRYLTRRRRIEIAHMLTLTERQIKIWFQNRRMKWKKKEHK
    AKNQISLLGSHK*NEKSF
11 >osedax_3_assembly_cap
12 LNVGPNSSQRRRGRQTYSTRYQTLELEKEFQFNHYLTRKRRIEIAHVLCCLTERQIKIWFQNRRMK
    LKKEKQQIKDL
```

```
13 >osedax_1_transcriptome_owe
14 PFYPWM-GVVGPNSSQRRRGRQTYSTRYQTLELEKEFQFNHYLTRKRRIEIAHVLCLTERQIKIW
   FQNRMMKLKKEKQQIKDL
15 >osedax_2_transcriptome_owe
16 PIFPWMRRMHLGDGIDGIDSKRTRTSYTRHQTLLELEKEFHFNRYLTRRRRIEIAHMLTLTERQIK
   IWFQNRMMKWKKEHK
17 >osedax_3_transcriptome_owe
18 IFPWMKKVHNGTSNGGFNGVDSKRSRTAYTRHQVLELEKEFHFNRYLTRRRRIEIAHTLCLTER
   QIKIWFQNRMMKWKDHK
19 >osedax_1_transcriptome_cap
20 GPNSNQRRRGRQTYTRYQTLELEKEFKYNRYLTRRRRIELSHTLCLTERQIKIWFQNRMMKEKK
   EIQAIEL
21 >osedax_2_transcriptome_cap
22 GPNSNQRRRGRQTYTRYQTLELEKEFKYNRYLTRRRRIELSHTLCLTERQIKIWFQNRMMKEKK
   EIQAIEL
23 >osedax_3_transcriptome_cap
24 GPNSNQRRRGRQTYTRYQTLELEKEFKYNRYLTRRRRIELSHTLCLTERQIKIWFQNRMMKEKK
   EIQAIEL
```

check hoxes found yesterday in osedax in its genome:

osedax.fa

```
1 >osedax_Lox2
2 ASEPSPNVMPWMSIVGPNSNQRRRGRQTYTRYQTLELEKEFKYNRYLTRRRRIELSHTLCLT
   ERQIKIWFQNR
3 >osedax_Hox1
4 VATYKWM TVKR NPKTVK QTPQSSDYNGNSSTTSGACCASGSHFRSSPLSPSHSPSSIGSGCGG
   GNLSGNGLPPNLGR TNFTNKQLTELEKEFHFNRYLTRARRIEIAASLCLNETQVKIWFQNRMMK
   QKKRLKEGHAAQWTV DTE
5 >osedax_Hox4
6 RTAYTRHQVLELEKEFHFNRYLTRRRRIEIAHTLCLTERQIKIWFQNRMMKWK
```

osedax_nucl.fa

used the previous sequences and translated them into nucleotide with this [online tool](#)

```
1 >osedax_Lox2
2 GCCAGCGAGCCCCCAGCCCCAACGTGATGTACCCCTGGATGAGCATCGTGGGCCCCAAC
3 AGCAACCAGAGGAGGAGGGGCAGGCAGACCTACACCAGGTACCAGACCCTGGAGCTGGAG
4 AAGGAGTTCAAGTACAACAGGTACCTGACCAGGAGGAGGAGGATCGAGCTGAGCCACACC
5 CTGTGCCTGACCGAGAGGCAGATCAAGATCTGGTTCCAGAACAGGAGG
6 >osedax_Hox1
7 GTGGCCACCTACAAGTGGATGACCGTGAAGAGGAACGCCCCAAGACCGTGAAGCAGACC
8 CCCCAGAGCAGCGACTACAACGGCAACAGCAGCACCACCAGCGGCGCCTGCTGCGCCAGC
9 GGCAGCCACTTCAGGAGCAGCCCCCTGAGCCCCAGCCACAGCCCCAGCAGCATCGGCAGC
10 GGCTGCGGCGGCGGCAACCTGAGCGGCAACGGCCTGCCCCCAACCTGGGCAGGACCAAC
11 TTCACCAACAAGCAGCTGACCGAGCTGGAGAAGGAGTTCCACTTCAACAGGTACCTGACC
12 AGGGCCAGGAGGATCGAGATCGCCGCCAGCCTGTGCCTGAACGAGACCCAGGTGAAGATC
13 TGGTTCCAGAACAGGAGGATGAAGCAGAAGAAGAGGCTGAAGGAGGGCCACGCCGCCAG
14 TGGACCGTGGACACCGAG
15 >osedax_Hox4
16 AGGACCGCCTACACCAGGCACCAGGTGCTGGAGCTGGAGAAGGAGTTCCACTTCAACAGG
17 TACCTGACCAGGAGGAGGAGGATCGAGATCGCCACACCCTGTGCCTGACCGAGAGGCAG
18 ATCAAGATCTGGTTCCAGAACAGGAGGATGAAGTGGAAGAAG
```

check_osedax_assembly.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/H0X_genes/further/assembly/check/osedax
3 #$ -o /data/scratch/btx654/H0X_genes/further/assembly/check/osedax
4 #$ -j y
5 #$ -pe smp 8
6 #$ -l h_vmem=10G
7 #$ -l h_rt=120:0:0
8 #$ -l highmem
9
```

```
10 module load blast+
11 makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax
   /haploidization/purge_dups/osedax_purged.fa -dbtype nucl -out os
   edax_assembly
12 tblastn -db osedax_assembly -query osedax.fa -out osedax_assembl
   y_tblastn_out -max_target_seqs 5 -evalue 1e-10 -num_threads 8 -o
   utfmt 6
13 tblastn -db osedax_assembly -query osedax.fa -out osedax_assembl
   y_tblastn_out.html -max_target_seqs 5 -evalue 1e-10 -num_threads
   8 -html
14
15 blastn -db osedax_assembly -query osedax_nucl.fa -out osedax_ass
   embly_blastn_out -max_target_seqs 5 -evalue 1e-10 -num_threads 8
   -outfmt 6
16 blastn -db osedax_assembly -query osedax_nucl.fa -out osedax_ass
   embly_blastn_out.html -max_target_seqs 5 -evalue 1e-10 -num_thre
   ads 8 -html
```

```
seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax/tr
   inity/osedax_*/*.Trinity.fasta list_transcript > list_transcrip
   t.fa
```

check_osedax_blastn.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/H0X_genes/further/assembly/check/ose
   dax
3 #$ -o /data/scratch/btx654/H0X_genes/further/assembly/check/osed
   ax
4 #$ -j y
5 #$ -pe smp 8
6 #$ -l h_vmem=10G
7 #$ -l h_rt=120:0:0
8 #$ -l highmem
```

```
9
10 module load blast+
11
12 blastn -db osedax_assembly -query list_transcript.fa -out osedax
    _assembly_blastn_out -max_target_seqs 5 -evaluate 1e-10 -num_threa
    ds 8 -outfmt 6
13 blastn -db osedax_assembly -query list_transcript.fa -out osedax
    _assembly_blastn_out.html -max_target_seqs 5 -evaluate 1e-10 -num_
    threads 8 -html
```

hox5.fa

```
1 >Owenia_fusiformis_Hox5
2 MSLYSLKSPAAYNSFMSDSGGGGHREFTHSENPYRAYTSGYPYTSHTAAPSGSTHQNGTPTDYS
    SFSNPATQRLIHPSYNREDSPTPVNNKPM PAATT SVITSTSPQDYSIKSRTTTEFATTKSSSQ
    ITDRDSA VDSPSPTPGSPVSPGQPNKDNDKYVKEEEDGSDREDRDGEGGADNPNIQIYPWM
    RRVHLGHDQNGAETKRTRTSYTRHQTLELEKEFHFNRYLTRRRRIEIAHSLNLTERQIKIWFQN
    RRMKWKKEHKLAHLAKSQAKMLDLALAQRAAEAKMHHHAAHGHTLHL
3
```

tblastn_osedax_assembly.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/H0X_genes/further/assembly/Hox1
3 #$ -o /data/scratch/btx654/H0X_genes/further/assembly/Hox1
4 #$ -j y
5 #$ -pe smp 8
6 #$ -l h_vmem=10G
7 #$ -l h_rt=120:0:0
8 #$ -l highmem
9
10 module load blast+
11 makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax
    /haploidization/purge_dups/osedax_purged.fa -dbtype nucl -out os
```

```
osedax_assembly
```

```
12 tblastn -db osedax_assembly -query hox1.fa -out osedax_assembly_
tblastn_out -max_target_seqs 5 -evalue 1e-10 -num_threads 8 -out
fmt 6
13 tblastn -db osedax_assembly -query hox1.fa -out osedax_assembly_
tblastn_out.html -max_target_seqs 5 -evalue 1e-10 -num_threads 8
-html
```

Hox result table

GENE	Osedax	Riftia	Oasisia	Paraescarpi a	Lamellibrac hia
HOX1	<p>signal found in the trinity transcripto me</p> <p>VATYKWM TVKRNAP KTVKQTP QSSDYNG NSSTTSG ACCASGS HFRSSPLS PSHSPSSI GSGCGGG NLSGNGL PPNLGRT NFTNKQLT ELEKEFHF NRYLTRAR RIEIAASLC</p>	<p>Rpac_RPAC G0000001 9294.1/66- 123</p>	<p>Oalv_OALV G0000001 9490.1/61- 118</p>	<p>Pech_nbis- mrna-1015 1/36-93</p>	<p>Lluy_FUN_0 32673- T1/90-147</p>

	<p>LNETQVKI WFQNRRM KQKKRLKE GHAAQWT VDTE</p> <p>TRINITY_D N22632_c0 _g1_i6</p> <p>TRINITY_D N16248_c0 _g1_i5 no traces in the assembly</p>				
HOX2	<p>Ofra_OFRA G0000001 1248.1/15 5-212</p>	<p>signal found in the trinity transcripto me</p> <p>YSSGAVIN MCGPVAP PPASGGG GLTTPGH PRRLRTAY TNTQLLEL EKEFHFNK YLCRPRRI EIAASLDL TERQVKV WFQNRR</p>	<p>Oalv_OALV G0000001 9775.1/14 0-197</p>	<p>Pech_nbis- mrna-1015 3/29-86</p>	<p>Lluy_FUN_0 32670- T1/134-19 1</p>

MKFKRQT			
QPKSSDG			
VAMPGDD			
DFGSPAID			
STTVSDD			
SHSPLGVS			
SVGDKDA			
PSGDTCG			
DDCDKTP			
GVKSGND			
RSDAADS			
VSLGQRS			
MTDVDDS			
AIKCEDMS			
RGKAPSV			
DASPSPTD			
AVPFSNPL			
ARCDSRV			
DAQREFN			
MSQHMG			
APLQPPVL			
THLTGSD			
GMQTVRN			
VCHPYVTP			
GNVDTSL			
ANPSRPM			
LHGPPSLP			
HSEQNAR			
TSFPPTGH			
TSSLHPAD			
VPARQAM			
RGPSLPQT			
IYPPTAGIH			
RLAAPHTK			

	HNSYLGS ATSRRHAP YIDISNMS ASDDRRM DNYFPSN GSDTCVG STSQYGV TPVYVRH DMQYARD HRPPQQQ SCYGDVT QQQNFTR NQDMLNV PFQNNVN CMPFSQP ATGDSHA YTNLAGC YSQPGMT NNNYYQG AAQYGAD TGSIGGCE TDYTSGY DGRYMTS HINSAETT PPDGDGV SSSFPSLS EFCQITNY NYL			
--	--	--	--	--

	>TRINITY_ DN39953_c 0_g1_i1			
--	-----------------------------------	--	--	--

		>TRINITY_ DN28744_c 0_g1_i1			
		no traces in the assembly			
HOX3	Ofra_OFRA G0000001 1246.1/10 9-166	Rpac_RPAC G0000001 9878.1/23 9-296	Oalv_OALV G0000001 9773.1/20 8-265	nbis_mrna_ 10152	Lluy_FUN_0 32669- T1/164-22 1
HOX4	signal found in the haploid assembly RTAYTRH QVLELEKE FHFNRYLT RRRRIEIAH TLCLTERQ IKIWFQNR RMKWKK tig0000404 6 arrow arr ow pilon pil on - 108:269	signal found in the haploid assembly AFVNTV* WRGCVSV SGSTGSF NGDNKRT RTAYTRH QVLELEKE FHFNRYLT RRRRIEIAH TLCLTERQ IKIWFQNR RMKWKK >tig000212 07_pech	Oalv_OALV G0000001 9768.1/18 1-238	Pech_nbis- mrna-1015 4/167-224	Lluy_FUN_0 05888- T1/165-22 2
HOX5	Ofra_OFRA	signal found in	signal found in	Pech_nbis- mrna-1015	Lluy_FUN_0

G00000000 993.1/119- 176	the trinity transcripto me	the transcripto me	6/22-79	38049- T1/12-69
	RRMHLGH DGVNGVE TKRTRTSY TRHQTLEL EKEFHFN YLTRRRRIE IAHMLNLT ERQIKIWF QNRRMK WKKEHKM AHLAKAQ AQKLETQ MHVGSAD MTRKS	RRMHLGH DGVNGVE TKRTRTSY TRHQTLEL EKEFHFN YLTRRRRIE IAHMLNLT ERQIKIWF QNRRMK WKKEHKM AHLAKAQ AQKLETQL HVG TADM TRKS		
	>TRINITY_ DN18618_c 0_g2_i2	TRINITY_D N28412_c0 _g1_i1		
	>TRINITY_ DN18618_c 0_g1_i1	TRINITY_D N15117_c0 _g3_i1		
	>TRINITY_ DN18618_c 0_g2_i1	TRINITY_D N9861_c0_ g1_i2		
	no traces in the assembly	TRINITY_D N9861_c0_ g1_i1		

			no traces in the assembly		
Lox2	<p>signal found in the trinity transcriptome</p> <p>ASEPPSPN VMYPWM SIVGPNSN QRRRGRQ TYTRYQTL ELEKEFKY NRYLTRRR RIELSHTL CLTERQIKI WFQNRR</p> <p>TRINITY_DN2605_c0_g1_i5</p> <p>tig00000951 arrow arrow pilon pilon</p>	<p>signal found in the trinity transcriptome</p> <p>PNSNQRR RGRQTYT RYQTLELE KEFKFNRY LTRRRRIEL SHMLCLT ERQIKIWF QNRRMKE KKEIQAIK ELNEKEKT KGTPTTV MPTAK</p> <p>>TRINITY_DN72223_c0_g1_i1</p> <p>no traces in the assembly</p>	<p>Oalv_OALV G0000001 9413.1/215-272</p>	<p>Pech_nbis-mrna-1015 8/6-63</p>	<p>FUN_03804 5-T1 ? very high blast hit with Oasisia but not clustering well</p>
Lox4	<p>Ofra_OFRA G0000001 1363.1/62-</p>	<p>signal found in the trinity transcriptome</p>	<p>Oalv_OALV G0000001 9412.1/17</p>	<p>Pech_nbis-mrna-1015 7/6-63</p>	<p>FUN_03804 6-T1 ? very high</p>

	119	me PFYPWMG VVGPNSS QRRRGRQ TYSRYQTL ELEKEFQF NHYLTRKR RIEIAHALC LTERQIKI WFQNRR >TRINITY_ DN726_c1_ g2_i3 no traces in the assembly	5-232		blast hit with Oasisia but not clustering well
Lox5	Ofra_OFRA G0000000 0992.1/10- 67	signal found in the trinity transcripto me GYEQKRT RQTYTRY QTLELEKE FHYNRYLT RRRRIEIAH ALGLSER QIKIWFQN RRMKWKK ENNLSKLT	Oalv_OALV G0000001 9872.1/18 8-245	Pech_nbis- mrna-1015 5/184-241	Lluy_FUN_0 38047- T1/185-24 2

		<p>GPNGNDQ PVESTNG SVE</p> <p>>TRINITY_ DN38983_c 0_g1_i1</p> <p>>TRINITY_ DN2094_c1 _g1_i2</p> <p>no traces in the assembly</p>			
Post1	Ofra_OFRA G0000000 9996.1/48- 112	<p>signal found in the haploid assembly</p> <p>ESDPFRVR KPRHVLA RPTPLRPR KKRKPYTK EQISDLEQ EYLDTTYI TRPKRTEI AKRLHLTE RQVKIWF QNRRMKE KKTNNKN VNLFEI</p>	<p>signal found in the haploid assembly</p> <p>EPDPFRAR KPRHILAR PTPLRPRK KRKPYTKE QISELEQE YLDTTYIT RPKRTEIA KRLHLTER QVKIWFQ NRRMKEK KTNNKNIN LYEV</p>	Pech_nbis- mrna-9378 /22-79	Lluy_FUN_0 30985- T1/70-127

		>tig000354 75_pech	>tig000018 29 - 619630:61 9899		
Post2	Ofra_OFRA G0000001 4621.1/33 2-389	MISSING checked the assembly and the transcripto me	Oalv_OALV G0000001 9415.1/17 2-229	Pech_nbis- mrna-1015 9/1-45	Lluy_FUN_0 38044- T1/172-22 9
Antp	No Antp detected both in the haploid assembly and in the transcripto me	No Antp detected both in the haploid assembly and in the transcripto me	No Antp detected both in the haploid assembly and in the transcripto me	No Antp detected in non redundant proteome and mRNA	No Antp detected in non redundant proteome and mRNA

Plots

```

1 library(tidyverse)
2 library(dplyr)
3 library(ggplot2)
4 library(data.table)
5 library(gplots)
6 library(pheatmap)

```

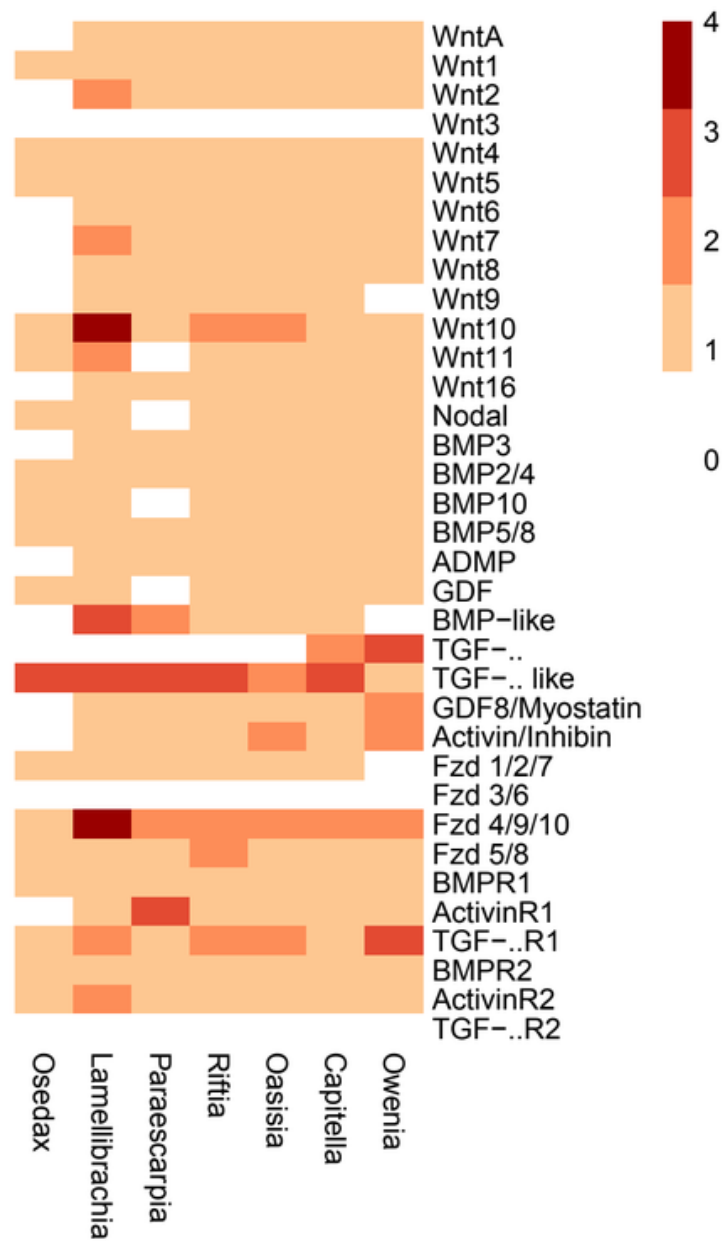


```

7 library(dendextend)
8 library(factoextra)
9 library(ComplexHeatmap)
10 library(RColorBrewer)
11 library(NbClust)
12 library(scales)
13
14 #Import data in matrix format
15 pathways <- read.delim("~/Desktop/pathways_updated.txt", row.names=1)
16
17 # Option 1. 0 to 1 relative abundance/expression (54 is the highest value in my dataset)
18
19 heatmap_color <- colorRampPalette(brewer.pal(n = 7, name = "OrRd"))(6)
20 heatmap_color[1] <- rgb(1,1,1)
21
22 # If you want common scale for different heatmaps:
23 # First define some "breaks"
24
25 pheatmap(pathways,
26           cluster_rows = FALSE,
27           cluster_cols = FALSE,
28           border_color = NA,
29           color = heatmap_color,
30           cellheight = 10,
31           cellwidth = 20)
32

```

- exported in PDF 8x6 inches



Summary plot

sort.sh

```
1 cut -f1 Developmental_pathways_table_unique_ok > header
2 cut -f2 Developmental_pathways_table_unique_ok > riftia
3 cut -f3 Developmental_pathways_table_unique_ok > oasisia
4 cut -f4 Developmental_pathways_table_unique_ok > osedax
5 cut -f5 Developmental_pathways_table_unique_ok > lamellibrachia
6 cut -f6 Developmental_pathways_table_unique_ok > owenia
7 cut -f7 Developmental_pathways_table_unique_ok > capitella
8 cut -f8 Developmental_pathways_table_unique_ok > paraescarpia
9
10 paste header osedax lamellibrachia paraescarpia riftia oasisia c
    apitella owenia > Developmental_pathways_table_unique_perfect
11 rm Developmental_pathways_table_unique_ok
12 rm header
13 rm osedax
14 rm lamellibrachia
15 rm paraescarpia
16 rm riftia
17 rm oasisia
18 rm capitella
19 rm owenia
```

Now in my folder I have only 7 files each containing a set of genes without header (eg Wnt genes)

sum.sh

```
1 input=$1
2 cut -f 2 $input | paste -sd+ - | bc > temp1
3 cut -f 3 $input | paste -sd+ - | bc > temp2
4 cut -f 4 $input | paste -sd+ - | bc > temp3
5 cut -f 5 $input | paste -sd+ - | bc > temp4
```

```

6 cut -f 6 $input | paste -sd+ - | bc > temp5
7 cut -f 7 $input | paste -sd+ - | bc > temp6
8 cut -f 8 $input | paste -sd+ - | bc > temp7
9 paste temp1 temp2 temp3 temp4 temp5 temp6 temp7 > "$input"_sum
10 rm temp*

```

tot number of genes per each group:

```

1      12 BMPlig.txt
2      6 BMPrec.txt
3      4 Fzd.txt
4      6 Ligands.txt
5      7 Receptors.txt
6      18 TranscriptionFactors.txt
7      13 Wnt.txt

```

```

1 library(tidyverse)
2 library(dplyr)
3 library(ggplot2)
4 library(data.table)
5 library(gplots)
6 library(pheatmap)
7 library(dendextend)
8 library(factoextra)
9 library(ComplexHeatmap)
10 library(RColorBrewer)
11 library(NbClust)
12 library(scales)
13
14 #Import data in matrix format
15 first_half_log <- read.delim("~/Desktop/Developmental_pathways/n

```

```
o isoforms/first_half_log_Paraescarpia_STAT=-inf", row.names=1)
16 second_half_log_signal <- read.delim("~/Desktop/Developmental_pa
    thways/no isoforms/second_half_log_signal", row.names=1)
17 second_half_log_receptor <- read.delim("~/Desktop/Developmental_
    pathways/no isoforms/second_half_log_receptor", row.names=1)
18 # Option 1. 0 to 1 relative abundance/expression (54 is the high
    est value in my dataset)
19
20 heatmap_color <- colorRampPalette(brewer.pal(n = 7, name = "Red
    s"))(1000)
21
22 # If you want common scale for different heatmaps:
23 # First define some "breaks"
24
25 pheatmap(first_half_log,
26           cluster_rows = FALSE,
27           cluster_cols = FALSE,
28           border_color = NA,
29           color = heatmap_color,
30           cellheight = 10,
31           cellwidth = 20)
32
33 pheatmap(second_half_log_signal,
34           cluster_rows = FALSE,
35           cluster_cols = FALSE,
36           border_color = NA,
37           color = heatmap_color,
38           cellheight = 10,
39           cellwidth = 20)
40
41 pheatmap(second_half_log_receptor,
```

```
42     cluster_rows = FALSE,  
43     cluster_cols = FALSE,  
44     border_color = NA,  
45     color = heatmap_color,  
46     cellheight = 10,  
47     cellwidth = 20)
```

```
1  library(tidyverse)  
2  library(dplyr)  
3  library(ggplot2)  
4  library(data.table)  
5  library(gplots)  
6  library(pheatmap)  
7  library(dendextend)  
8  library(factoextra)  
9  library(ComplexHeatmap)  
10 library(RColorBrewer)  
11 library(NbClust)  
12 library(scales)  
13  
14 #Import data in matrix format  
15 A <- read.delim("/Users/giacomo/Dropbox/11-Siboglinids/05-PAPER/  
00-DATA/04-Losses/03-DevelopmentalPathways/Summary/Wnt.txt_sum",  
row.names=1)  
16 B <- read.delim("/Users/giacomo/Dropbox/11-Siboglinids/05-PAPER/  
00-DATA/04-Losses/03-DevelopmentalPathways/Summary/Fzd.txt_sum",  
row.names=1)  
17 C <- read.delim("/Users/giacomo/Dropbox/11-Siboglinids/05-PAPER/  
00-DATA/04-Losses/03-DevelopmentalPathways/Summary/BMPlig.txt_su  
m", row.names=1)  
18 D <- read.delim("/Users/giacomo/Dropbox/11-Siboglinids/05-PAPER/  
00-DATA/04-Losses/03-DevelopmentalPathways/Summary/BMPrec.txt_su
```

```
m", row.names=1)
19 E <- read.delim("/Users/giacomo/Dropbox/11-Siboglinids/05-PAPER/
00-DATA/04-Losses/03-DevelopmentalPathways/Summary/Transcription
Factors.txt_sum", row.names=1)
20 F <- read.delim("/Users/giacomo/Dropbox/11-Siboglinids/05-PAPER/
00-DATA/04-Losses/03-DevelopmentalPathways/Summary/Ligands.txt_s
um", row.names=1)
21 G <- read.delim("/Users/giacomo/Dropbox/11-Siboglinids/05-PAPER/
00-DATA/04-Losses/03-DevelopmentalPathways/Summary/Receptors.txt
_sum", row.names=1)
22
23 # Option 1. 0 to 1 relative abundance/expression (54 is the high
est value in my dataset)
24 rescale_custom <- function(x) (x/18)
25 A_normalised <- t(apply(A, 1, rescale_custom))
26 rescale_custom <- function(x) (x/6)
27 B_normalised <- t(apply(B, 1, rescale_custom))
28 rescale_custom <- function(x) (x/14)
29 C_normalised <- t(apply(C, 1, rescale_custom))
30 rescale_custom <- function(x) (x/7)
31 D_normalised <- t(apply(D, 1, rescale_custom))
32 rescale_custom <- function(x) (x/920)
33 E_normalised <- t(apply(E, 1, rescale_custom))
34 rescale_custom <- function(x) (x/50)
35 F_normalised <- t(apply(F, 1, rescale_custom))
36 rescale_custom <- function(x) (x/66)
37 G_normalised <- t(apply(G, 1, rescale_custom))
38
39 # To make 0 a different colour
40 # First create whatever gradient (e.g. RdBu)
41 heatmap_color <- colorRampPalette(brewer.pal(n = 7, name = "Red
s"))(1000)
```

```
42 heatmap_color[1] <- rgb(1,1,1)
43 #column_labels = c("your","labels"),
44 #row_labels = c("your","labels"))
45
46 paletteLength <- 1000
47 # to go from 0 to max.value (e.g. 1):
48 myBreaks <- c(seq(1/paletteLength, 1, length.out=floor(paletteLength)))
49
50 pheatmap(A_normalised,
51         cluster_rows = FALSE,
52         cluster_cols = FALSE,
53         border_color = NA,
54         color = heatmap_color,
55         height = 25,
56         width = 20,
57         breaks = myBreaks)
58
59 pheatmap(B_normalised,
60         cluster_rows = FALSE,
61         cluster_cols = FALSE,
62         border_color = NA,
63         color = heatmap_color,
64         height = 25,
65         width = 20,
66         breaks = myBreaks)
67
68 pheatmap(C_normalised,
69         cluster_rows = FALSE,
70         cluster_cols = FALSE,
```



```

71     border_color = NA,
72     color = heatmap_color,
73     height = 25,
74     width = 20,
75     breaks = myBreaks)
76
77 pheatmap(D_normalised,
78     cluster_rows = FALSE,
79     cluster_cols = FALSE,
80     border_color = NA,
81     color = heatmap_color,
82     height = 25,
83     width = 20,
84     breaks = myBreaks)
85
86 pheatmap(E_normalised,
87     cluster_rows = FALSE,
88     cluster_cols = FALSE,
89     border_color = NA,
90     color = heatmap_color,
91     height = 25,
92     width = 20,
93     breaks = myBreaks)
94
95 pheatmap(F_normalised,
96     cluster_rows = FALSE,
97     cluster_cols = FALSE,
98     border_color = NA,
99     color = heatmap_color,
100    height = 25,

```

```
101     width = 20,  
102     breaks = myBreaks)  
103  
104 pheatmap(G_normalised,  
105     cluster_rows = FALSE,  
106     cluster_cols = FALSE,  
107     border_color = NA,  
108     color = heatmap_color,  
109     height = 25,  
110     width = 20,  
111     breaks = myBreaks)
```

1.txt

```
1 grouptot_genes  
2 Wnt13  
3 BMP1ig12  
4 BMPrec6  
5 Fzd4
```

2.txt

```
1 grouptot_genes  
2 TranscriptionFactors18  
3 Ligands6  
4 Receptors7
```

```
1 library(tidyverse)  
2 library(dplyr)  
3 library(ggplot2)  
4 library(data.table)  
5 library(gplots)  
6 library(pheatmap)
```

```
7 library(dendextend)
8 library(factoextra)
9 library(ComplexHeatmap)
10 library(RColorBrewer)
11 library(NbClust)
12 library(scales)
13
14 #Import data in matrix format
15 total_1 <- read.delim("/Users/giacomo/Dropbox/11-Siboglinids/05-
PAPER/00-DATA/04-Losses/03-DevelopmentalPathways/Summary/1.txt",
row.names=1)
16 total_2 <- read.delim("/Users/giacomo/Dropbox/11-Siboglinids/05-
PAPER/00-DATA/04-Losses/03-DevelopmentalPathways/Summary/2.txt",
row.names=1)
17
18 # Option 1. 0 to 1 relative abundance/expression (54 is the high
est value in my dataset)
19 rescale_custom <- function(x) (x/13)
20 total_1_normalised <- t(apply(total_1, 1, rescale_custom))
21 rescale_custom <- function(x) (x/18)
22 total_2_normalised <- t(apply(total_2, 1, rescale_custom))
23
24 # To make 0 a different colour
25 # First create whatever gradient (e.g. RdBu)
26 heatmap_color <- colorRampPalette(brewer.pal(n = 7, name = "Red
s"))(1000)
27 heatmap_color[1] <- rgb(1,1,1)
28 #column_labels = c("your","labels"),
29 #row_labels = c("your","labels"))
30
31 paletteLength <- 1000
```

```
32 # to go from 0 to max.value (e.g. 1):
33 myBreaks <- c(seq(1/paletteLength, 1, length.out=floor(paletteLength)))
34
35 pheatmap(total_1_normalised,
36           cluster_rows = FALSE,
37           cluster_cols = FALSE,
38           border_color = NA,
39           color = heatmap_color,
40           height = 25,
41           width = 20,
42           breaks = myBreaks)
43
44 pheatmap(total_2_normalised,
45           cluster_rows = FALSE,
46           cluster_cols = FALSE,
47           border_color = NA,
48           color = heatmap_color,
49           height = 25,
50           width = 20,
51           breaks = myBreaks)
```

Mega sum

summing everything together and then plotting bubbleplots

count_not_missing_genes.sh

```
1 input=$1
2 cut -f 2 $input | grep -cwv "0" > temp1
3 cut -f 3 $input | grep -cwv "0" > temp2
4 cut -f 4 $input | grep -cwv "0" > temp3
5 cut -f 5 $input | grep -cwv "0" > temp4
```

```
6 cut -f 6 $input | grep -cwv "0" > temp5
7 cut -f 7 $input | grep -cwv "0" > temp6
8 cut -f 8 $input | grep -cwv "0" > temp7
9 paste temp1 temp2 temp3 temp4 temp5 temp6 temp7 > "$input"_not_m
  issing_genes
10 rm temp*
```

```
1 library(tidyverse)
2 library(dplyr)
3 library(ggplot2)
4 library(data.table)
5 library(gplots)
6 library(pheatmap)
7 library(dendextend)
8 library(factoextra)
9 library(ComplexHeatmap)
10 library(RColorBrewer)
11 library(NbClust)
12 library(scales)
13
14 #Import data in matrix format
15 Averages <- read.delim("/Users/giacomo/Dropbox/11-Siboglinids/05
  -PAPER/00-DATA/04-Losses/03-DevelopmentalPathways/Summary/mega_s
  um/all_together.txt_sum", row.names=1)
16
17 # Option 1. 0 to 1 relative abundance/expression (54 is the high
  est value in my dataset)
18 rescale_custom <- function(x) (x/1032)
19 Averages_normalised <- t(apply(Averages, 1, rescale_custom))
20
```

```

21 # To make 0 a different colour
22 # First create whatever gradient (e.g. RdBu)
23 heatmap_color <- colorRampPalette(brewer.pal(n = 7, name = "Red
  s"))(1000)
24 heatmap_color[1] <- rgb(1,1,1)
25 #column_labels = c("your","labels"),
26 #row_labels = c("your","labels"))
27
28 paletteLength <- 1000
29 # to go from 0 to max.value (e.g. 1):
30 myBreaks <- c(seq(1/paletteLength, 1, length.out=floor(paletteLe
  ngth)))
31
32 pheatmap(Averages_normalised,
33           cluster_rows = FALSE,
34           cluster_cols = FALSE,
35           border_color = NA,
36           color = heatmap_color,
37           height = 25,
38           width = 20,
39           breaks = myBreaks)

```

Chapter 5

Matrix MetalloProteinases

Intro

Gene	Name	Aliases	Location	Description
------	------	---------	----------	-------------

MMP1	Interstitial collagenase	CLG, CLGN	secreted	Substrates include Col I, II, III, VII, VIII, X, gelatin
MMP2	Gelatinase-A, 72 kDa gelatinase		secreted	Substrates include Gelatin, Col I, II, III, IV, Vii, X
MMP3	Stromelysin 1	CHDS6, MMP-3, SL-1, STMY, STMY1, STR1	secreted	Substrates include Col II, IV, IX, X, XI, gelatin
MMP7	Matrilysin, PUMP 1	MMP-7, MPSL1, PUMP-1	secreted	membrane associated through binding to cholesterol sulfate in cell membranes, substrates include: fibronectin, laminin, Col IV, gelatin
MMP8	Neutrophil collagenase	CLG1, HNC, MMP-8, PMNL-CL	secreted	Substrates include Col I, II, III, VII, VIII, X, aggrecan, gelatin
MMP9	Gelatinase-B, 92 kDa gelatinase	CLG4B, GELB, MANDP2, MMP-9	secreted	Substrates include Gelatin, Col IV, V

MMP10	Stromelysin 2	SL-2, STMY2	secreted	Substrates include Col IV, laminin, fibronectin, elastin
MMP11	Stromelysin 3	SL-3, ST3, STMY3	secreted	MMP-11 shows more similarity to the MT-MMPs, is convertase-activatable and is secreted therefore usually associated to convertase-activatable MMPs. Substrates include Col IV, fibronectin, laminin, aggrecan
MMP12	Macrophage metalloelastase	HME, ME, MME, MMP-12	secreted	Substrates include elastin, fibronectin, Col IV
MMP13	Collagenase 3	CLG3, MANDP1, MMP-13	secreted	Substrates include Col I, II, III, IV, IX, X, XIV, gelatin

MMP14	MT1-MMP	MMP-14, MMP-X1, MT- MMP, MT- MMP 1, MT1- MMP, MT1MMP, MTMMP1, WNCHRS	membrane- associated	type-I transmembra ne MMP; substrates include gelatin, fibronectin, laminin
MMP15	MT2-MMP	MT2-MMP, MTMMP2, SMCP-2, MMP-15, MT2MMP	membrane- associated	type-I transmembra ne MMP; substrates include gelatin, fibronectin, laminin
MMP16	MT3-MMP	C8orf57, MMP-X2, MT- MMP2, MT- MMP3, MT3- MMP	membrane- associated	type-I transmembra ne MMP; substrates include gelatin, fibronectin, laminin
MMP17	MT4-MMP	MT4-MMP, MMP-17, MT4MMP, MTMMP4	membrane- associated	glycosyl phosphatidyl inositol- attached; substrates include fibrinogen, fibrin
<i>MMP18</i>	Collagenase 4, xcol4, xenopus		–	No known human

	collagenase			orthologue
MMP19	RASI-1, occasionally referred to as stromelysin-4	MMP18, RASI-1, CODA	–	
MMP20	Enamelysin	AI2A2, MMP-20	secreted	
MMP21	X-MMP	MMP-21, HTX7	secreted	Our findings suggest that MMP-21 functions in embryogenesi s and tumor progression.
MMP23A	CA-MMP		membrane- associated	type-II transmembra ne cysteine array
MMP23B	–	MIFR, MIFR-1, MMP22, MMP23A	membrane- associated	type-II transmembra ne cysteine array
MMP24	MT5-MMP	MMP-24, MMP25, MT- MMP 5, MT- MMP5, MT5- MMP, MT5MMP, MTMMP5	membrane- associated	type-I transmembra ne MMP
MMP25	MT6-MMP	MMP-25, MMP20, MMP20A, MMPL1, MT- MMP 6, MT-	membrane- associated	glycosyl phosphatidyl inositol- attached

		MMP6, MT6-MMP, MT6MMP, MTMMP6		
MMP26	Matrilysin-2, endometase		–	
MMP27	MMP-22, C-MMP	MMP-27	–	
MMP28	Epilysin	EPILYSIN, MM28, MMP-25, MMP-28, MMP25	secreted	Discovered in 2001 and given its name due to have been discovered in human keratinocytes . Unlike other MMPs this enzyme is constitutively expressed in many tissues (Highly expressed in testis and at lower levels in lung , heart , brain , colon , intestine , placenta , salivary glands , uterus , skin). A threonine

				replaces proline in its cysteine switch (PRCGVTD). ^[1 4]
--	--	--	--	---

Focus

Col I is

The gained collagenases related GO terms in Osedax are:

- 1 GO:0030574 collagen catabolic process (way more present in osedax)
- 2 GO:0032963 collagen metabolic process (not too much more present in osedax)

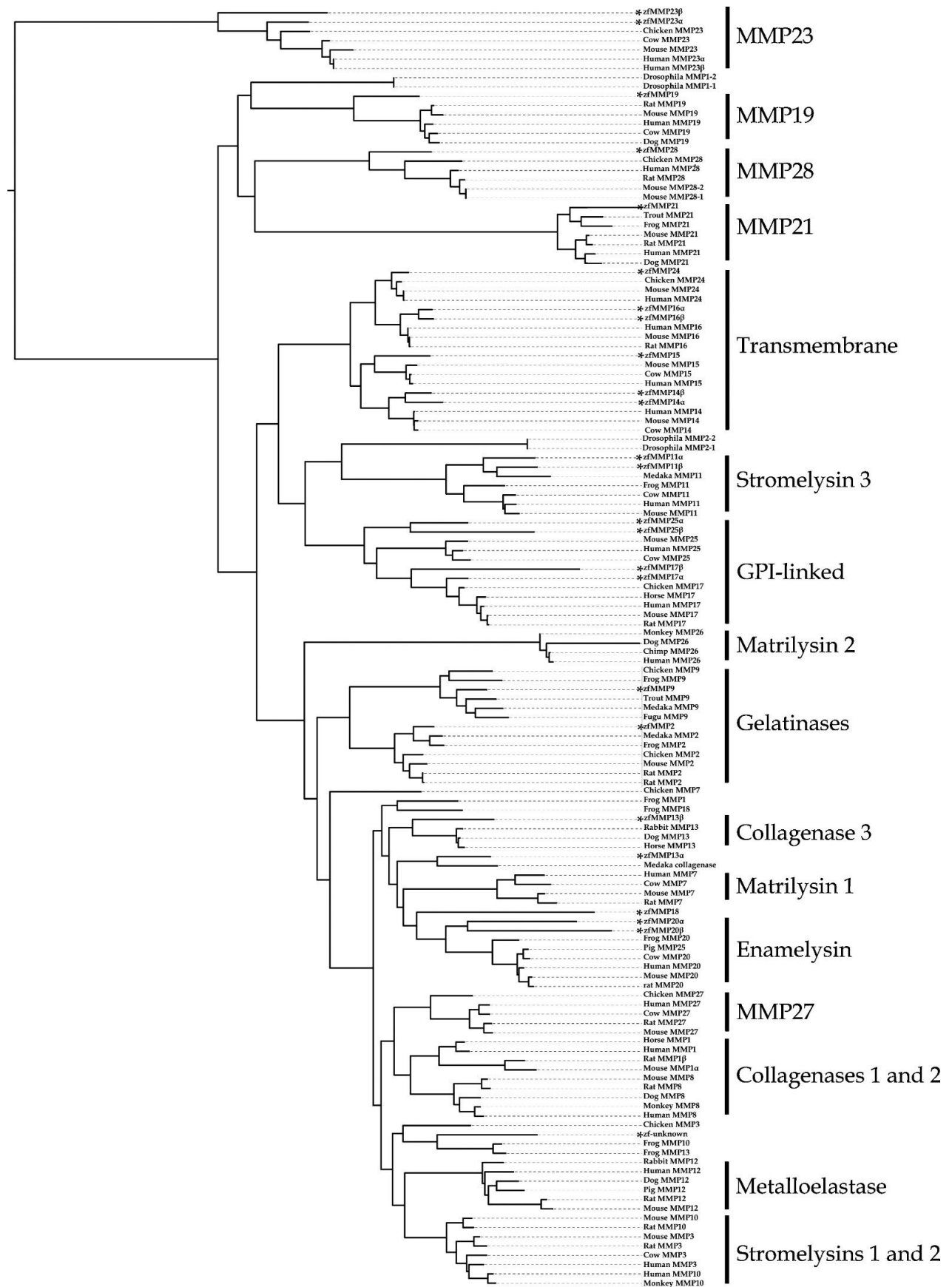
Collagen related Panther are:

- 1 PTHR10201 MATRIX METALLOPROTEINASE (way more present in osedax. even more than GOterm)
- 2 PTHR12411 CYSTEINE PROTEASE FAMILY C1-RELATED (not too much more present. cathepsin)

From [this](#) article I will be able to download the Matrix Metallo Proteases (MMP) sequences of Human, Mouse and Zebrafish: [Link to additional material](#)

I think I will have to manually collect the sequences from [here](#)

I can use this as a template



- MMP23
- MMP19
- MMP28
- MMP21
- MMP24
- MMP16
- MMP15
- MMP14
- MMP11
- MMP25
- MMP17
- MMP26
- MMP9
- MMP2
- MMP13
- MMP7
- MMP20
- MMP27
- MMP1
- MMP8
- MMP12
- MMP10
- MMP3
- MMP18 (one from osedax japonicus) NOT PRESENT IN HUMAN

CODE

test.sh

```
1  #!/bin/bash
2
3  species=osedax
4  xls="$species"_isoform.xls
5  output="$species"_MMPs.txt
6  output2="$species"_MMPs_count.txt
7
```

```

8  cut -f 18,19,20 /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL
   _isoforms_annotations/$xls | fgrep "PTHR10201" | cut -f 1,2,3 >
   $output
9  cut -f 18,19,20 /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL
   _isoforms_annotations/$xls | fgrep "PTHR10201" | wc -l > $output
   2
10
11
12 species=oasisia
13 xls="$species"_isoform.xls
14 output="$species"_MMPs.txt
15 output2="$species"_MMPs_count.txt
16
17 cut -f 18,19,20 /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL
   _isoforms_annotations/$xls | fgrep "PTHR10201" | cut -f 1,2,3 >
   $output
18 cut -f 18,19,20 /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL
   _isoforms_annotations/$xls | fgrep "PTHR10201" | wc -l > $output
   2
19
20
21 species=riftia
22 xls="$species"_isoform.xls
23 output="$species"_MMPs.txt
24 output2="$species"_MMPs_count.txt
25
26 cut -f 18,19,20 /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL
   _isoforms_annotations/$xls | fgrep "PTHR10201" | cut -f 1,2,3 >
   $output
27 cut -f 18,19,20 /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL
   _isoforms_annotations/$xls | fgrep "PTHR10201" | wc -l > $output
   2

```

```
28
29
30 species=lamellibrachia
31 xls="$species"_isoform.xls
32 output="$species"_MMPs.txt
33 output2="$species"_MMPs_count.txt
34
35 cut -f 18,19,20 /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL
_isoforms_annotations/$xls | fgrep "PTHR10201" | cut -f 1,2,3 >
$output
36 cut -f 18,19,20 /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL
_isoforms_annotations/$xls | fgrep "PTHR10201" | wc -l > $output
2
37
38
39 species=paraescarpia
40 xls="$species"_isoform.xls
41 output="$species"_MMPs.txt
42 output2="$species"_MMPs_count.txt
43
44 cut -f 18,19,20 /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL
_isoforms_annotations/$xls | fgrep "PTHR10201" | cut -f 1,2,3 >
$output
45 cut -f 18,19,20 /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL
_isoforms_annotations/$xls | fgrep "PTHR10201" | wc -l > $output
2
46
47
48 species=owenia
49 xls="$species"_isoform.xls
50 output="$species"_MMPs.txt
51 output2="$species"_MMPs_count.txt
```



```
52
53 cut -f 2,3,10 /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL_i
soforms_annotaions/$xls | fgrep "PTHR10201" | cut -f 1,2 > $out
put
54 cut -f 2,3,10 /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL_i
soforms_annotaions/$xls | fgrep "PTHR10201" | wc -l > $output2
55
56 species=capitella
57 xls="$species"_isoform.xls
58 output="$species"_MMPs.txt
59 output2="$species"_MMPs_count.txt
60
61 cut -f 1,2,3 /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL_is
oforms_annotaions/$xls | fgrep "PTHR10201" | cut -f 1,2 > $outp
ut
62 cut -f 1,2,3 /data/SBCS-MartinDuranLab/03-Giacomo/data/00-ALL_is
oforms_annotaions/$xls | fgrep "PTHR10201" | wc -l > $output2
```

seqtk.sh

```
1 #!/bin/bash
2 grep "PTHR10201" /data/SBCS-MartinDuranLab/03-Giacomo/data/00-AL
L_isoforms_annotaions/capitella* | cut -f1 | sed 's/^/Ctel_/' >
MMPs_candidates_capitella
3 module load seqtk
4 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/C
tel.fa MMPs_candidates_capitella > MMPs_candidates_capitella.fa
5 grep "PTHR10201" /data/SBCS-MartinDuranLab/03-Giacomo/data/00-AL
L_isoforms_annotaions/owenia*.xls | cut -f2 > MMPs_candidates_
owenia
6 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/O
fus.fa MMPs_candidates_owenia > MMPs_candidates_owenia.fa
7 grep "PTHR10201" /data/SBCS-MartinDuranLab/03-Giacomo/data/00-AL
L_isoforms_annotaions/riftia* | cut -f2 | sed 's/^/Rpac_/' > MM
```

```
Ps_candidates_riftia
8 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/R
  pac.fa MMPs_candidates_riftia > MMPs_candidates_riftia.fa
9 grep "PTHR10201" /data/SBCS-MartinDuranLab/03-Giacomo/data/00-AL
  L_isoforms_annotations/oasisia* | cut -f2 | sed 's/^/Oalv_/' > M
  MMPs_candidates_oasisia
10 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/O
  alv.fa MMPs_candidates_oasisia > MMPs_candidates_oasisia.fa
11 grep "PTHR10201" /data/SBCS-MartinDuranLab/03-Giacomo/data/00-AL
  L_isoforms_annotations/osedax* | cut -f2 | sed 's/^/Ofra_/' > MM
  Ps_candidates_osedax
12 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/O
  fra.fa MMPs_candidates_osedax > MMPs_candidates_osedax.fa
13 grep "PTHR10201" /data/SBCS-MartinDuranLab/03-Giacomo/data/00-AL
  L_isoforms_annotations/lamellibrachia* | cut -f2 | sed 's/^/Lluy
  _/' > MMPs_candidates_lamellibrachia
14 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/L
  luy.fa MMPs_candidates_lamellibrachia > MMPs_candidates_lamellib
  rachia.fa
15 grep "PTHR10201" /data/SBCS-MartinDuranLab/03-Giacomo/data/00-AL
  L_isoforms_annotations/paraescarpia* | cut -f2 | sed 's/nbis_mrn
  a_/nbis-mrna-/' | sed 's/^/Pech_/' > MMPs_candidates_paraescarp
  ia
16 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/P
  ech.fa MMPs_candidates_paraescarpia > MMPs_candidates_paraescarp
  ia.fa
```

```
1 bash seqtk.sh
2 cat MMPs_candidates*.fa > MMPs_candidates_all.fa
```

Phylogenetic Tree

Step 1 - MMPs phylogenetic tree

make an alignment using the txt file obtained before and the online tool [MAFFT](#)

don't select any additional options

export in fasta format and name the output:

```
MMPs_sequences_MAFFT1.fasta
```

Step 2 - MMPs phylogenetic tree

download and install [Jalview](#)

open the file:

```
MMPs_sequences_MAFFT1.fasta
```

Cut away all the areas with no conservation. Basically we should maintain only the domain

and save the cutted file as:

```
MMPs_sequences_Jalview.fasta
```

Step 3 - MMPs phylogenetic tree

make an alignment again using the file "MMPs_sequences_Jalview.fasta" obtained before and the online tool [MAFFT](#)

select the option "L-INS-i"

export in fasta format and name the output:

```
MMPs_sequences_MAFFT2.fasta
```

Step 4 - MMPs phylogenetic tree

use Jalview again. Open the file:

```
MMPs_sequences_MAFFT2.fasta
```

No chopping this time, just open the file and export it to fasta as:

```
MMPs_sequences_Jalview2.fasta
```

Step 5 - MMPs phylogenetic tree

```
MMPs_sequences_gblocks.fa
```

or use trimAl

on my personal pc

```
1 conda create -n fasttree_env
2 conda activate fasttree_env
3 conda install -c bioconda fasttree
4 conda install -c bioconda trimal
```

```
trimal -in 4a.fasta -out MMPs_sequences_trimal.fasta
```

Step 6 - MMPs phylogenetic tree

make a tree

```
1 conda activate fasttree_env
2 FastTree MMPs_sequences_trimal.fasta > MMPs_sequences_a.tree
```

and then open it with FigTree

Improved

I need to use human MMPs to find the gene families corresponding to those genes in all the species I have used for Orthofinder

Step1 - Blast

First I need to blast the MMPs sequences against the Human non redundant proteome to find the correct gene name

test.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/MMPs/blast
3 #$ -o /data/scratch/btx654/MMPs/blast
4 #$ -j y
5 #$ -pe smp 8
6 #$ -l h_vmem=10G
7 #$ -l h_rt=120:0:0
8 #$ -l highmem
9
10 $non_redundant_prot=/data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/Hsap.fasta
11 cp /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/Hsap.fasta ./
12
13 #make a diamond BLAST database of this proteome and BLAST the consensus.fa.classified dataset against it, to find potential bona fide genes. To make sure we only get the real genes, the e-value
```

is very stringent.

14 `module load anaconda3`

15 `source activate diamond`

16 `diamond makedb --in $Hsap.fa -d Hsap_non_redundant_prot`

17

18 `diamond blastp -d Hsap_non_redundant_prot -q MMPs_Hsap_genes.txt`
`-o default.1e10.blastp -f 6 qseqid bitscore evalue stitle -k 25`
`-e 1e-10 -p 8`

19 `diamond blastp -d Hsap_non_redundant_prot -q MMPs_Hsap_genes.txt`
`-o ultra_sensitive.1e10.blastp --ultra-sensitive -f 6 qseqid bit`
`score evalue stitle -k 25 -e 1e-10 -p 8`

		OG000007 1	OG003260 3	OG000756 0	OG000110 5
MMP23-	Hsap_ENS T00000356 026		0		
MMP19-	Hsap_ENS T00000322 569	0			
MMP28-	Hsap_ENS T00000605 424	0			
MMP21-	Hsap_ENS T00000368 808			0	

MMP24-	Hsap_ENS T00000246 186	0			
MMP16-	Hsap_ENS T00000286 614	0			
MMP15-	Hsap_ENS T00000219 271	0			
MMP14-	Hsap_ENS T00000311 852	0			
MMP11-	Hsap_ENS T00000215 743	0			
MMP25-	Hsap_ENS T00000336 577	0			
MMP17-	Hsap_ENS T00000360 564	0			
MMP26-	Hsap_ENS T00000380 390	0			

MMP9-	Hsap_ENS T00000372 330				0
MMP2-	Hsap_ENS T00000219 070				0
MMP13-	Hsap_ENS T00000340 273	0			
MMP7-	Hsap_ENS T00000260 227	0			
MMP20-	Hsap_ENS T00000260 228	0			
MMP27-	Hsap_ENS T00000260 229	0			
MMP1-	Hsap_ENS T00000315 274	0			
MMP8-	Hsap_ENS T00000236 826	0			

MMP12-	Hsap_ENS T00000571 244	0			
MMP10-	Hsap_ENS T00000279 441	0			
MMP3-	Hsap_ENS T00000299 855	0			

OG0000071

Hsap|ENST00000215743, Hsap|ENST00000219271, Hsap|ENST00000236826, Hsap|ENST00000246186, Hsap|ENST00000260227, Hsap|ENST00000260228, Hsap|ENST00000260229, Hsap|ENST00000279441, Hsap|ENST00000286614, Hsap|ENST00000299855, Hsap|ENST00000311852, Hsap|ENST00000315274, Hsap|ENST00000322569, Hsap|ENST00000336577, Hsap|ENST00000340273, Hsap|ENST00000360564, Hsap|ENST00000380390, Hsap|ENST00000571244, Hsap|ENST00000605424

OG0000071	275	28	Eumetazoa	0
Hrob,Ofra				

OG0001105

Hsap|ENST00000219070, Hsap|ENST00000339841, Hsap|ENST00000344839, Hsap|ENST00000372330

Step2 - get the sequences

test.sh

```
1 module load seqtk
2
```



```
3 head -1 /data/SBCS-MartinDuranLab/03-Giacomo/data/all_together/gene_family_evolution/orthofinder_Jun2021/ultra_sensitive/Results_Jun09/Orthogroups/Orthogroups.csv | cut -f 2-29 | sed 's/\t\t*/\n/g' > species_list.txt
4
5 grep "OG0000071\|OG0032603\|OG0007560\|OG0001105" /data/SBCS-MartinDuranLab/03-Giacomo/data/all_together/gene_family_evolution/orthofinder_Jun2021/ultra_sensitive/Results_Jun09/Orthogroups/Orthogroups.csv | cut -f 2-29 > MMPs_Orthogroups.csv
6
7 for run in {1..28}; do
8 species=$(cat species_list.txt | head -"$run" | tail -1)
9 echo $species
10
11 cut -f "$run" MMPs_Orthogroups.csv > "$species"_MMPs_Orthogroups.csv
12 sed 's/|/_/g' "$species"_MMPs_Orthogroups.csv | sed 's/, /\n/g' | sed -r '/^\s*$/d' > "$species"_gene_names.txt
13 rm "$species"_MMPs_Orthogroups.csv
14
15 seqtk subseq /data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/"$species".fa "$species"_gene_names.txt > "$species"_genes.fa
16
17 done
```

```
bash test.sh
```

Step3 - Tree

Collect all the sequences extracted in step2 and add them to those downloaded from uniprot (except those from 5 sibo + owenia and capi and those from human)

```
cat Blan_genes.fa Bpla_genes.fa Cgig_genes.fa Dgyr_genes.fa Eand_genes.fa Gaeg_genes.fa Hmia_genes.fa Hrob_genes.fa Lana_genes.f
```

```
a Lgig_genes.fa Locu_genes.fa Myes_genes.fa Ngen_genes.fa Nvec_g
enes.fa Paus_genes.fa Skow_genes.fa Smar_genes.fa Smed_genes.fa
Spur_genes.fa Tcas_genes.fa > ALL_genes.fa
```

Then I am adding these sequences to “2.fasta” obtained after aligning and trimming only the MMPs sequences downloaded from uniprot. I DID NOT ADD THE SEQUENCES FROM SIBO AND CAPI-OWENIA OBTAINED IN THE SECTION CODE

The tree is still a bit messy:

I am going to name some of the clades that came up in the tree:

MMP21_like

MMP11_17_25_like

MMP

MMPA (is the branch containing osedax japonicus)

MMPB

MMPC

MMPD

MMPE

MMPF

MMPG

removed:

```
1 >OFUSG03511.1
2 >Pech_nbis-mrna-5872
3 >Lluy_FUN_007880-T1
```

raxml.sh

```
1 #!/bin/bash
2 #$ -pe smp 5
3 #$ -l highmem
4 #$ -l h_vmem=10G
5 #$ -l h_rt=240:0:0
6 #$ -cwd
```

```
7  # $ -j y
8  module load raxml
9  raxmlHPC -f a -b 476 -p 903 -x 12345 -# autoMRE -m PROTGAMMAAUTO
   -s MMPs_sequences_trimal_ok.fa -n MMPs_raxml_tree.tre
```

removed: (after cropping with Jalview there were no ammino acids left)

```
1  Nvec_ED029979_MMP2_9_like
2  Nvec_ED033047_MMP2_9_like
3  Nvec_ED032111_MMP2_9_like
```

I am now aligning just the ZnMc domain which contains the active site:
after trimal some proteins result to not have that domain so I will remove them:

```
1  >Ofra_OFrag000000002650.1
2  >Lluy_FUN_026149-T1_MMP2_9_like
3  >Blan_BL22912_cuf1_MMP21_like
4  >Ofra_OFrag000000014433.1
5  >Nvec_ED045372_MMP2_9_like
6  Nvec_ED045377_MMP2_9_like
7  >Paus_g3879.t1_MMP2_9_like
8  >Ofra_OFrag00000000712.1
9  >Ofra_OFrag00000000713.1
10 >Ofra_OFrag00000000714.1
11 >Ofra_OFrag00000000715.1
12 >Ofra_OFrag00000000716.1
13 >Ofra_OFrag00000000717.1
14 >Ofra_OFrag00000000718.1
15 >Ofra_OFrag00000000719.1
16 >Ofra_OFrag00000000720.1
17 >Ofra_OFrag00000000721.1
18 >Ofra_OFrag00000000722.1
```

```
19 >Ofra_OFRAG000000014432.1
20 >Nvec_ED029980_MMPB
21 >Ngen_g12182.t1_MMP?
22 >Lana_g4594.t1_MMP2_9_like
23 >Ngen_g7295.t1_MMP2_9_like
24 >Ngen_g37516.t1_MMP2_9_like
25 >Ofra_OFRAG000000014430.1
26 >Ofra_OFRAG000000014431.1
27 >Locu_ENSLOCT000000011484_MMP2_9_like
28 >Ngen_g22040.t1_MMP2_9_like
29 >Nvec_ED047418_MMP2_9_like
30 >Paus_g15030.t1_MMP2_9_like
31 >Ngen_g21330.t1_MMP2_9_like
32 >Blan_BL96786_cuf2_MMP2_9_like
33 >Blan_BL10256_cuf0_MMP2_9_like
34 >Paus_g10720.t1_MMP2_9_like
35 >Ofra_OFRAG000000002369.1
36 >Ngen_g23379.t1_MMP2_9_like
37 >Nvec_ED035188_MMP?
38 >Nvec_ED045821_MMP2_9_like
39 >Ofra_OFRAG000000002286.1
40 >Ngen_g34577.t1_MMPB
41 >Ofra_OFRAG000000007211.1
42 >Ofra_OFRAG000000007212.1
43 >Ofra_OFRAG000000007213.1
44 >Ofra_OFRAG000000008179.1
45 >Blan_BL10324_cuf0_MMP2_9_like
46 >Ofra_OFRAG000000013465.1
47
48 Nvec_ED029979_MMP2_9_like
```

```
49 Nvec_ED033047_MMP2_9_like
50 Nvec_ED032111_MMP2_9_like
51
52 removed from MMPs_sequences_trimal_ok.fa:
53 >Lana_g26508.t1_MMP2_9_like
54 >Ngen_g7296.t1_MMP2_9_like
55 >Locu_ENSLOCT000000017075_MMP23
56 >Blan_BL06531_evm0_MMP?
57 >OFUSG11076.1_MMP?
```

Raxml-ng

on the cluster

```
1 module load anaconda3
2 conda create --prefix /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/raxml_env
3 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/raxml_env
4 conda install -c bioconda raxml
  • raxml-8.2.12
```

raxml.sh

```
1 #!/bin/bash
2 #$ -pe smp 7
3 #$ -l highmem
4 #$ -l h_vmem=20G
5 #$ -l h_rt=240:0:0
6 #$ -cwd
7 #$ -j y
8
9 module load anaconda3
10 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/raxml_env
```

```
11
12 raxmlHPC -f a -b 476 -p 903 -# autoMRE -m PROTGAMMAAUTO -s MMPs_
sequences_trimal_ok.fa -n MMPs_raxml_tree.tre
```

raxml_oceane.sh

```
1 #!/bin/bash
2 #$ -pe smp 7
3 #$ -l highmem
4 #$ -l h_vmem=20G
5 #$ -l h_rt=240:0:0
6 #$ -cwd
7 #$ -j y
8
9 module load anaconda3
10 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
3/raxml_env
11
12 raxmlHPC -f a -b 476 -p 903 -x 12345 -# autoMRE -m PROTGAMMAAUTO
-s MMPs_sequences_trimal_ok.fa -n MMPs_raxml_tree_oceane.tre
```

```
1 module load anaconda3
2 conda create --prefix /data/SBCS-MartinDuranLab/03-Giacomo/src/a
naconda3/raxml_ng_env
3 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
3/raxml_ng_env
4 conda install -c bioconda raxml-ng
```

on command line:

```
1 raxml-ng --check --msa MMPs_sequences_trimal_ok.fa --model LG+G
8+F --prefix T1
2 raxml-ng --parse --msa T1.raxml.reduced.phy --model LG+G8+F --pr
efix T2
```

raxml_ng.sh

```
1  #!/bin/bash
2  #$ -pe smp 20
3  #$ -l highmem
4  #$ -l h_vmem=20G
5  #$ -l h_rt=240:0:0
6  #$ -cwd
7  #$ -j y
8
9  module load anaconda3
10 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
    3/raxml_ng_env
11
12 raxml-ng --all --msa T2.raxml.rba --model LG+G8+F --bs-trees 100
    0 --threads 20 --prefix MMPs_raxml_ng
```

- LG+G8+F Perform an all-in-one analysis (ML tree search + non-parametric bootstrap) (10 randomized parsimony starting trees, fixed empirical substitution matrix (LG), empirical aminoacid frequencies from alignment, 8 discrete GAMMA categories, 1000 bootstrap replicates):

```
1  module load anaconda3
2  conda create --prefix /data/SBCS-MartinDuranLab/03-Giacomo/src/a
    naconda3/modeltest_env
3  conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
    3/modeltest_env
4  conda install -c bioconda modeltest-ng
```

modeltest_ng.sh

```
1  #!/bin/bash
2  #$ -pe smp 20
3  #$ -l highmem
4  #$ -l h_vmem=20G
```

```
5  #$ -l h_rt=240:0:0
6  #$ -cwd
7  #$ -j y
8
9  module load anaconda3
10 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
    3/modeltest_env
11
12 /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/modeltest_env
    /bin/modeltest-ng -i MMPs_sequences_trimal_ok.fa --datatype aa
    --processes 20 --template raxml -t ml --output MMPs_modeltest_ng
    _output
```

IQtree

```
1  module load anaconda3
2  conda create --prefix /data/SBCS-MartinDuranLab/03-Giacomo/src/a
    naconda3/IQtree_env
3  conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
    3/IQtree_env
4  conda install -c bioconda iqtree
```

iqtree.sh

```
1  #!/bin/bash
2  #$ -pe smp 20
3  #$ -l highmem
4  #$ -l h_vmem=20G
5  #$ -l h_rt=240:0:0
6  #$ -cwd
7  #$ -j y
8
9  module load anaconda3
```



```
10 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
    3/IQtree_env
11
12 iqtree -s MMPs_sequences_trimal_ok.fas -m MFP -B 1000
```

Blast_pacbio

in order to confirm the duplication events I am going to blast the ZnMc domain of Osedax against the raw pacbio reads

```
1 module load seqtk
2 seqtk seq -a osedax_pb_raw.fastq > osedax_pb_raw.fa
```

ZnMc_domain_osedax.fa

```
1 >ZnMc_domain_osedax
2 GFVWKHLNITYKITEYTRKVSHTHIDEAAAKALNFWGEVTQLNFRQVSPYSKADIDIKFVVGDH
  GDGLPFDGKGKIIGHAFPPEYGLAHFDDAESWVIGECDDASINILQVMTHEFGHSLGLAHSFNR
  SNVMFPSYKGYEPNFALSGDDIKGIQSLYG
```

blast.sh

```
1 #!/bin/bash
2 #$ -cwd
3 #$ -j y
4 #$ -pe smp 8
5 #$ -l h_vmem=20G
6 #$ -l h_rt=120:0:0
7 #$ -l highmem
8
9 module load blast+
10 makeblastdb -in osedax_pb_raw.fa -dbtype nucl -out osedax_pacbio
11 tblastn -db osedax_pacbio -query ZnMc_domain_osedax.fa -out osed
  ax_pacbio_tblastn_out -max_target_seqs 5 -evalue 1e-10 -num_thre
  ads 8 -outfmt 6
12 tblastn -db osedax_pacbio -query ZnMc_domain_osedax.fa -out osed
  ax_pacbio_tblastn_out.html -max_target_seqs 5 -evalue 1e-10 -num
```

```
_threads 8 -html
```

used the previous sequences and translated them into nucleotide with this [online tool](#)

ZnMc_domain_osedax_nucl.fa

```
1 >ZnMc_domain_osedax
2 GGCTTCGTGTGGAAGCACCTGAACATCACCTACAAGATCACCGAGTACACCCGCAAGGTG
3 AGCCACACCCACATCGATGAGGCCGCCGCCAAGGCCCTGAACTTCTGGGGCGAGGTGACC
4 CAGCTGAACTTCCGCCAGGTGAGCCCCTACAGCAAGGCCGATATCGATATCAAGTTCGTG
5 GTGGGCGATCACGGCGATGGCCTGCCCTTCGATGGCAAGGGCGGCATCATCGGCCACGCC
6 TTCCCCCCCCGAGTACGGCCTGGCCCACTTCGATGATGCCGAGAGCTGGGTGATCGGCGAG
7 TGCGATGATGCCAGCATCAACATCCTGCAGGTGATGACCCACGAGTTCGGCCACAGCCTG
8 GGCCTGGCCACAGCTTCAACCGCAGCAACGTGATGTTCCCAGCTACAAGGGCTACGAG
9 CCCAACTTCGCCCTGAGCGGCGATGATATCAAGGGCATCCAGAGCCTGTACGGC
```

blast_nucl.sh

```
1 #!/bin/bash
2 #$ -cwd
3 #$ -j y
4 #$ -pe smp 8
5 #$ -l h_vmem=20G
6 #$ -l h_rt=120:0:0
7 #$ -l highmem
8
9 module load blast+
10 blastn -db osedax_pacbio -query ZnMc_domain_osedax_nucl.fa -out
    osedax_pacbio_blastn_out -max_target_seqs 5 -evalue 1e-10 -num_t
    hreads 8 -outfmt 6
11 blastn -db osedax_pacbio -query ZnMc_domain_osedax_nucl.fa -out
    osedax_pacbio_blastn_out.html -max_target_seqs 5 -evalue 1e-10 -
    num_threads 8 -html
```

no hits found

blast_assembly.sh

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -pe smp 8
5  #$ -l h_vmem=20G
6  #$ -l h_rt=120:0:0
7  #$ -l highmem
8
9  module load blast+
10 makeblastdb -in /data/SBCS-MartinDuranLab/03-Giacomo/data/osedax
    /haploidization/purge_dups/osedax_purged.fa -dbtype nucl -out os
    edax_assembly
11 tblastn -db osedax_assembly -query ZnMc_domain_osedax.fa -out os
    edax_pacbio_tblastn_out -max_target_seqs 5 -evaluate 1e-10 -num_th
    reads 8 -outfmt 6
12 tblastn -db osedax_assembly -query ZnMc_domain_osedax.fa -out os
    edax_pacbio_tblastn_out.html -max_target_seqs 5 -evaluate 1e-10 -n
    um_threads 8 -html
```

```
1  module load anaconda3
2  conda create -n genomePolishing python=2.7 anaconda
3  conda activate genomePolishing
4  conda install -c bioconda genomicconsensus
5  conda install -c bioconda pbmm2
6  conda install -c bioconda pbgcpp
7  conda install pbbam
```

First Round:

Align raw PB reads to the assembly in a bam file

Code to run pbmm2:

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step
   1
3  #$ -j y
4  #$ -o /data/SBCS-MartinDuranLab/03-Giacomo/logs/genomePolishing
5  #$ -pe smp 5
6  #$ -l h_vmem=8G
7  #$ -l h_rt=72:0:0
8
9  cd /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step1
10
11  module load anaconda2
12  source activate genomePolishing
13
14  pbmm2 align \
15  /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/canu/Oasisia.
   contigs.fasta \
16  /data/SBCS-MartinDuranLab/03-Giacomo/data/oasisia/00-pacbio/dat
   a2/pb/r64044_20190812_215729/1_A01/m64044_190812_220643.subread
   s.bam \
17  /data/scratch/btx604/Oasisia/genomePolishing/pbalign/step1/oasi
   sia_pbalign_step1.bam
18
19  module load samtools/1.9
20
21  samtools sort /data/scratch/btx604/Oasisia/genomePolishing/pbali
   gn/step1/oasisia_pbalign_step1.bam -o /data/scratch/btx604/Oasis
   ia/genomePolishing/pbalign/step1/oasisia_sorted_step1.bam
22  samtools index /data/scratch/btx604/Oasisia/genomePolishing/pbal
   ign/step1/oasisia_sorted_step1.bam
```

- Sorted output can be generated using `--sort`

code to run minimap2

minimap2.sh

```
1 #!/bin/bash
2 #$ -cwd
3 #$ -j y
4 #$ -pe smp 10
5 #$ -l h_vmem=8G
6 #$ -l h_rt=120:0:0
7
8 module load minimap2
9 minimap2 -ax map-pb /data/SBCS-MartinDuranLab/03-Giacomo/data/os
  edax/haploidization/purge_dups/osedax_purged.fa /data/scratch/bt
  x654/MMPs/blast_pacbio/osedax_pb_raw.fastq > raw_pacbio_vs_purge
  d_assembly.sam
```

```
1 #!/bin/bash
2 #$ -cwd
3 #$ -j y
4 #$ -pe smp 5
5 #$ -l h_vmem=4G
6 #$ -l h_rt=120:0:0
7 module load samtools
8 samtools view -S -b raw_pacbio_vs_purged_assembly.sam > raw_pacb
  io_vs_purged_assembly.bam
9 samtools sort raw_pacbio_vs_purged_assembly.bam -o raw_pacbio_vs
  _purged_assembly_sorted.bam
10 samtools index raw_pacbio_vs_purged_assembly_sorted.bam
11 samtools view raw_pacbio_vs_purged_assembly.bam "tig00006601|arr
  ow|arrow|pilon|pilon" > tig00006601.bam
```

R studio

```
1 library(tidyverse)
2 library(dplyr)
3 library(ggplot2)
4 library(data.table)
5 library(gplots)
6 library(pheatmap)
7 library(dendextend)
8 library(factoextra)
9 library(ComplexHeatmap)
10 library(RColorBrewer)
11 library(NbClust)
12 library(scales)
13
14 #Import data in matrix format
15 df1 <- read.delim("~/Dropbox/02-OweniaGenome/02-Figures/00-Clutter/02-Fran&YanFigures/Urechis_Hox_genes_16_stages.txt")
16 df2
17 #Name columns and stuff
18 names_row <- raw_data[,1]
19 df_curated <- data.matrix(raw_data[,2:ncol(raw_data)])
20 rownames(df_curated) <- names_row
21
22 # Option 1. Deviation from the mean (z-score).
23 # Normalise your data around the mean (0) and then look deviation around mean.
24 df_normalised <- t(scale(t(df_curated)))
25 df_normalised <- na.omit(df_normalised)
26
27 # Option 1. 0 to 1 relative abundance/expression
28 rescale_custom <- function(x) (x/(max(x)))
```

```
29 df_normalised <- t(apply(df_curated, 1, rescale_custom))
30 df_normalised <- na.omit(df_normalised)
31
32 # Colours
33 # Option 1. RdBu gradient.
34 heatmap_color <- colorRampPalette(rev(brewer.pal(n = 7, name = "
RdBu")))(100)
35
36 # Option 2. Reds gradient.
37 heatmap_color <- colorRampPalette(rev(brewer.pal(n = 7, name = "
Reds")))(100)
38
39 # Option 3. 1/2 RdBu gradient.
40 heatmap_color <- colorRampPalette(rev(brewer.pal(n = 7, name = "
RdBu")))(100)
41 heatmap_color <- heatmap_color[50:100]
42
43 # Option 4. RdGy gradient.
44 heatmap_color <- colorRampPalette(rev(brewer.pal(n = 7, name = "
RdGy")))(100)
45
46 # To make 0 a different colour
47 # First create whatever gradient (e.g. RdBu)
48 heatmap_color <- colorRampPalette(rev(brewer.pal(n = 7, name = "
RdBu")))(100)
49 heatmap_color[1] <- rgb(0,0,0) # here include the colour you're
interested in
50
51 # Conversion from dataframe to datamatrix (if necessary)
52 dm <- data.matrix(df)
53
```

```
54 # Heatmaps
55 # ComplexHeatmap package
56 ComplexHeatmap::Heatmap(dm,
57                           cluster_columns = FALSE,
58                           cluster_rows = FALSE,
59                           col = heatmap_color,
60                           heatmap_legend_param = list(color_bar =
"continuous"))
61
62 #column_labels = c("your","labels"),
63 #row_labels = c("your","labels"))
64
65 # pheatmap
66 pheatmap(df,
67           cluster_rows = FALSE,
68           cluster_cols = FALSE,
69           border_color = NA,
70           color = heatmap_color)
71
72 # If you want common scale for different heatmaps:
73 # First define some "breaks"
74
75 paletteLength <- 100
76 # to go from 0 to max.value (e.g. 1):
77 myBreaks <- c(seq(max.value/paletteLength, max.value, length.out
=floor(paletteLength)))
78 # e.g. if it was actually 1, then:
79 myBreaks <- c(seq(1/paletteLength, 1, length.out=floor(paletteLe
ngth)))
80
81 # Then use the following:
```



```
82 pheatmap(df,
83     cluster_rows = FALSE,
84     cluster_cols = FALSE,
85     border_color = NA,
86     color = heatmap_color,
87     breaks = myBreaks)
88
89 # To set sizes of heatmaps:
90 # Option 1: Cell height, each cell is the same size
91 pheatmap(df,
92     cluster_rows = FALSE,
93     cluster_cols = FALSE,
94     border_color = NA,
95     color = heatmap_color,
96     cellheight = 20,
97     cellwidth = 20)
98
99 # Option 2: Heatmap height, the whole heatmap will have this size
100 pheatmap(df,
101     cluster_rows = FALSE,
102     cluster_cols = FALSE,
103     border_color = NA,
104     color = heatmap_color,
105     height = 20,
106     width = 20)
```

Expression

```
1 sed 's/^.>0fra_//g' osedax_IDs_only > osedax_IDs_only_ok
2 grep -f osedax_IDs_only_ok osedax_kallisto_tpm.tsv > osedax_MMPs_expression
```

```
3 for S in $(cat osedax_IDs_only_ok | awk '{print $1}'); do grep
  $$ osedax_MMPs_expression; done > right_order_osedax_MMPs_expres
  sion
```

```
1 #scale all the gene values then subset
2 #in R
3 A <- read.delim("/Users/giacomo/Dropbox/11-Siboglinids/05-PAPER/
  00-DATA/05-Expanssions/01-Collagenases/RNAseq/osedax_kallisto_tp
  m.tsv", row.names=1)
4 B <- data.matrix(scale(A))
5 write.csv(B, file="/Users/giacomo/Dropbox/11-Siboglinids/05-PAPE
  R/00-DATA/05-Expanssions/01-Collagenases/RNAseq/osedax_kallisto_
  tpm_scaled.tsv")
6 #in shell
7 sed 's/\\,/\\t/g' osedax_kallisto_tpm_scaled.tsv | sed 's/\\\\"//g' >
  osedax_kallisto_tpm_scaled_ok.tsv
8 cut -f 1,2,3 osedax_kallisto_tpm_scaled_ok.tsv > osedax_kallisto
  _tpm_scaled_ok_ok.tsv
9 grep -f osedax_IDs_only_ok osedax_kallisto_tpm_scaled_ok_ok.tsv
  > osedax_MMPs_expression_ALLscaled
10 for S in $(cat osedax_IDs_only_ok | awk '{print $1}'); do grep
  $$ osedax_MMPs_expression_ALLscaled; done > right_order_osedax_M
  MPs_expression_ALLscaled
```

```
1 library(tidyverse)
2 library(dplyr)
3 library(ggplot2)
4 library(data.table)
5 library(gplots)
6 library(pheatmap)
7 library(dendextend)
8 library(factoextra)
```

```
9 library(ComplexHeatmap)
10 library(RColorBrewer)
11 library(NbClust)
12 library(scales)
13
14 #Import data in matrix format
15 A <- read.delim("/Users/giacomo/Dropbox/11-Siboglinids/05-PAPER/
00-DATA/05-Expanssions/01-Collagenases/RNAseq/right_order_osedax
_MMPs_expression", row.names=1)
16 A_scaled <- data.matrix(scale(A))
17 # Option 1. 0 to 1 relative abundance/expression (54 is the high
est value in my dataset)
18 B <- read.delim("/Users/giacomo/Dropbox/11-Siboglinids/05-PAPER/
00-DATA/05-Expanssions/01-Collagenases/RNAseq/right_order_osedax
_MMPs_expression_ALLscaled", row.names=1)
19
20 # To make 0 a different colour
21 # First create whatever gradient (e.g. RdBu)
22 heatmap_color <- colorRampPalette(brewer.pal(n = 7, name = "Red
s"))(1000)
23 heatmap_color[1] <- rgb(1,1,1)
24 #column_labels = c("your","labels"),
25 #row_labels = c("your","labels"))
26
27 paletteLength <- 1000
28 # to go from 0 to max.value (e.g. 1):
29 myBreaks <- c(seq(1/paletteLength, 1, length.out=floor(paletteLe
ngth)))
30
31 pheatmap(A_scaled,
32           cluster_rows = FALSE,
33           cluster_cols = FALSE,
```

```
34     border_color = NA,  
35     color = heatmap_color,  
36     height = 5,  
37     width = 25)  
38  
39 pheatmap(B,  
40     cluster_rows = FALSE,  
41     cluster_cols = FALSE,  
42     border_color = NA,  
43     color = heatmap_color,  
44     height = 5,  
45     width = 25)
```

B is definitely the best one!

Panther Metabolism

Nomenclature:

A_II_1 (compound_pathwayNumber_stepNumber)

A_II_1A (uppercase final letters in case more than one enzyme can catalise the reaction "OR")

A_II_1a (lowercase final letters in case more than one enzyme are needed to catalise the reaction "AND")

A_II_1_3 (in case the same enzyme can catalise different steps)

CODE

```
1  #!/bin/bash  
2  
3  species=$1  
4  xls="$species"*xls  
5  output="$species"_metabolism.txt  
6  
7  cut -f 18,19 $xls > "$species"_annotations
```

```
8 annotations="$species"_annotations
9
10 while read line; do
11     Panther_ID=$(cut -f 2 <<< "$line")
12     echo $Panther_ID
13     genes=$(fgrep "$Panther_ID" $annotations | cut -f 1)
14     echo $genes
15     echo $genes >> $output
16 done < metabolism_pantherID.txt
```

test_owenia_capitella.sh

```
1 #!/bin/bash
2
3 species=owenia
4 xls=Owenia_annotation_v250920.1_TrinoPantherK0.xls
5 output="$species"_metabolism.txt
6 #cut -f 18,19 for the 5 siboglinidae, -f 1,2 for capitella
7 cut -f 2,10 $xls > "$species"_annotations
8 annotations="$species"_annotations
9
10 while read line; do
11     Panther_ID=$(cut -f 2 <<< "$line")
12     echo $Panther_ID
13     genes=$(fgrep "$Panther_ID" $annotations | cut -f 1)
14     echo $genes
15     echo $genes >> $output
16 done < metabolism_pantherID.txt
17
18 species=capitella
19 xls=Capitella_annotation_Feb2021_TrinoPantherK0.xls
20 output="$species"_metabolism.txt
```

```
21 #cut -f 18,19 for the 5 siboglinidae, -f 1,2 for capitella
22 cut -f 1,2 $xls > "$species"_annotations
23 annotations="$species"_annotations
24
25 while read line; do
26     Panther_ID=$(cut -f 2 <<< "$line")
27     echo $Panther_ID
28     genes=$(fgrep "$Panther_ID" $annotations | cut -f 1)
29     echo $genes
30     echo $genes >> $output
31 done < metabolism_pantherID.txt
```

Glycine degradation

GlycineDegradation_pantherID.txt

```
1 GD_1      PTHR46120
2 GD_2      PTHR13847:SF187
3 GD_3      PTHR13847:SF200
4 GD_4      PTHR11680
```

script.sh

```
1 #!/bin/bash
2
3 species=osedax
4 xls="$species"_isoform.xls
5 output="$species"_GlycineDegradation.txt
6 output2="$species"_GlycineDegradation_count.txt
7
8 while read line; do
9     Panther_ID=$(cut -f 2 <<< "$line")
```

```
10     genes=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | cut -f 18)
11     echo $genes >> $output
12     count=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | wc -l)
13     echo $count >> $output2
14 done < GlycineDegradation_pantherID.txt
15
16 species=oasisia
17 xls="$species"_isoform.xls
18 output="$species"_GlycineDegradation.txt
19 output2="$species"_GlycineDegradation_count.txt
20
21 while read line; do
22     Panther_ID=$(cut -f 2 <<< "$line")
23     genes=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | cut -f 18)
24     echo $genes >> $output
25     count=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | wc -l)
26     echo $count >> $output2
27 done < GlycineDegradation_pantherID.txt
28
29 species=riftia
30 xls="$species"_isoform.xls
31 output="$species"_GlycineDegradation.txt
32 output2="$species"_GlycineDegradation_count.txt
33
34 while read line; do
35     Panther_ID=$(cut -f 2 <<< "$line")
36     genes=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
```

```
omo/data/00-ALL_isoforms_annotations/$xls | cut -f 18)
37     echo $genes >> $output
38     count=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | wc -l)
39     echo $count >> $output2
40 done < GlycineDegradation_pantherID.txt
41
42 species=lamellibrachia
43 xls="$species"_isoform.xls
44 output="$species"_GlycineDegradation.txt
45 output2="$species"_GlycineDegradation_count.txt
46
47 while read line; do
48     Panther_ID=$(cut -f 2 <<< "$line")
49     genes=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | cut -f 18)
50     echo $genes >> $output
51     count=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | wc -l)
52     echo $count >> $output2
53 done < GlycineDegradation_pantherID.txt
54
55 species=paraescarpia
56 xls="$species"_isoform.xls
57 output="$species"_GlycineDegradation.txt
58 output2="$species"_GlycineDegradation_count.txt
59
60 while read line; do
61     Panther_ID=$(cut -f 2 <<< "$line")
62     genes=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | cut -f 18)
```



```
63     echo $genes >> $output
64     count=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | wc -l)
65     echo $count >> $output2
66 done < GlycineDegradation_pantherID.txt
67
68 species=owenia
69 xls="$species"_isoform.xls
70 output="$species"_GlycineDegradation.txt
71 output2="$species"_GlycineDegradation_count.txt
72
73 while read line; do
74     Panther_ID=$(cut -f 2 <<< "$line")
75     genes=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | cut -f 2)
76     echo $genes >> $output
77     count=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | wc -l)
78     echo $count >> $output2
79 done < GlycineDegradation_pantherID.txt
80
81 species=capitella
82 xls="$species"_isoform.xls
83 output="$species"_GlycineDegradation.txt
84 output2="$species"_GlycineDegradation_count.txt
85
86 while read line; do
87     Panther_ID=$(cut -f 2 <<< "$line")
88     genes=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | cut -f 1)
89     echo $genes >> $output
```

```

90     count=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | wc -l)
91     echo $count >> $output2
92 done < GlycineDegradation_pantherID.txt

```

```

1 cut -f 1 GlycineDegradation_pantherID.txt > GlycineDegradation_p
antherID_FirstColumn
2 paste GlycineDegradation_pantherID_FirstColumn owenia_GlycineDeg
radation_count.txt capitella_GlycineDegradation_count.txt lamell
ibrachia_GlycineDegradation_count.txt paraescarpia_GlycineDegrad
ation_count.txt oasisia_GlycineDegradation_count.txt riftia_Glyc
ineDegradation_count.txt osedax_GlycineDegradation_count.txt > G
lycineDegradation_count_table.txt

```

Glycine Cleavage

GlycineCleavageSystem_pantherID.txt

```

1 GCS_1      PTHR22912:SF151
2 GCS_2      PTHR11773:SF1
3 GCS_3      PTHR43757:SF2

```

script.sh

```

1 #!/bin/bash
2
3 species=osedax
4 xls="$species"_isoform.xls
5 output="$species"_GlycineCleavageSystem.txt
6 output2="$species"_GlycineCleavageSystem_count.txt
7
8 while read line; do
9     Panther_ID=$(cut -f 2 <<< "$line")
10    genes=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac

```

```
omo/data/00-ALL_isoforms_annotations/$xls | cut -f 18)
11     echo $genes >> $output
12     count=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | wc -l)
13     echo $count >> $output2
14 done < GlycineCleavageSystem_pantherID.txt
15
16 species=oasisia
17 xls="$species"_isoform.xls
18 output="$species"_GlycineCleavageSystem.txt
19 output2="$species"_GlycineCleavageSystem_count.txt
20
21 while read line; do
22     Panther_ID=$(cut -f 2 <<< "$line")
23     genes=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | cut -f 18)
24     echo $genes >> $output
25     count=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | wc -l)
26     echo $count >> $output2
27 done < GlycineCleavageSystem_pantherID.txt
28
29 species=riftia
30 xls="$species"_isoform.xls
31 output="$species"_GlycineCleavageSystem.txt
32 output2="$species"_GlycineCleavageSystem_count.txt
33
34 while read line; do
35     Panther_ID=$(cut -f 2 <<< "$line")
36     genes=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | cut -f 18)
```

```
37     echo $genes >> $output
38     count=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | wc -l)
39     echo $count >> $output2
40 done < GlycineCleavageSystem_pantherID.txt
41
42 species=lamellibrachia
43 xls="$species"_isoform.xls
44 output="$species"_GlycineCleavageSystem.txt
45 output2="$species"_GlycineCleavageSystem_count.txt
46
47 while read line; do
48     Panther_ID=$(cut -f 2 <<< "$line")
49     genes=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | cut -f 18)
50     echo $genes >> $output
51     count=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | wc -l)
52     echo $count >> $output2
53 done < GlycineCleavageSystem_pantherID.txt
54
55 species=paraescarpia
56 xls="$species"_isoform.xls
57 output="$species"_GlycineCleavageSystem.txt
58 output2="$species"_GlycineCleavageSystem_count.txt
59
60 while read line; do
61     Panther_ID=$(cut -f 2 <<< "$line")
62     genes=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | cut -f 18)
63     echo $genes >> $output
```

```
64     count=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | wc -l)
65     echo $count >> $output2
66 done < GlycineCleavageSystem_pantherID.txt
67
68 species=owenia
69 xls="$species"_isoform.xls
70 output="$species"_GlycineCleavageSystem.txt
71 output2="$species"_GlycineCleavageSystem_count.txt
72
73 while read line; do
74     Panther_ID=$(cut -f 2 <<< "$line")
75     genes=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | cut -f 2)
76     echo $genes >> $output
77     count=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | wc -l)
78     echo $count >> $output2
79 done < GlycineCleavageSystem_pantherID.txt
80
81 species=capitella
82 xls="$species"_isoform.xls
83 output="$species"_GlycineCleavageSystem.txt
84 output2="$species"_GlycineCleavageSystem_count.txt
85
86 while read line; do
87     Panther_ID=$(cut -f 2 <<< "$line")
88     genes=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
omo/data/00-ALL_isoforms_annotations/$xls | cut -f 1)
89     echo $genes >> $output
90     count=$(fgrep "$Panther_ID" /data/SBCS-MartinDuranLab/03-Giac
```

```
omo/data/00-ALL_isoforms_annotations/${xls} | wc -l)
```

```
91     echo $count >> $output2
```

```
92 done < GlycineCleavageSystem_pantherID.txt
```

```
1 cut -f 1 GlycineCleavageSystem_pantherID.txt > GlycineCleavageSystem_pantherID_FirstColumn
```

```
2 paste GlycineCleavageSystem_pantherID_FirstColumn owenia_GlycineCleavageSystem_count.txt capitella_GlycineCleavageSystem_count.txt lamellibrachia_GlycineCleavageSystem_count.txt paraescarpia_GlycineCleavageSystem_count.txt oasisia_GlycineCleavageSystem_count.txt riftia_GlycineCleavageSystem_count.txt osedax_GlycineCleavageSystem_count.txt > GlycineCleavageSystem_count_table.txt
```

```
3 echo "step"'\t'"owenia"'\t'"capitella"'\t'"lamellibrachia"'\t'"paraescarpia"'\t'"oasisia"'\t'"riftia"'\t'"osedax" > header.txt
```

```
4 cat header.txt GlycineCleavageSystem_count_table.txt > GlycineCleavageSystem_count_table_OK.txt
```

Visualisation

```
1 library(tidyverse)
```

```
2 library(dplyr)
```

```
3 library(ggplot2)
```

```
4 library(data.table)
```

```
5 library(gplots)
```

```
6 library(pheatmap)
```

```
7 library(dendextend)
```

```
8 library(factoextra)
```

```
9 library(ComplexHeatmap)
```

```
10 library(RColorBrewer)
```

```
11 library(NbClust)
```

```
12 library(scales)
```

```
13
```

```
14 A <- read.delim("/Users/giacomo/Desktop/GlycineDegradation/Glyci
neDegradation_count_table.txt ", row.names=1)
15
16 heatmap_color <- colorRampPalette(brewer.pal(n = 7, name = "Red
s"))(3)
17 #column_labels = c("your","labels"),
18 #row_labels = c("your","labels"))
19
20 paletteLength <- 3
21 pheatmap(A,
22         cluster_rows = FALSE,
23         cluster_cols = FALSE,
24         border_color = NA,
25         color = heatmap_color,
26         height = 25,
27         width = 20)
```

Expression

```
cd /Users/giacomo/Dropbox/11-Siboglinids/05-PAPER/00-DATA/06-Met
abolism/01-AA/expression
```

```
1 while read line; do
2   echo $line
3   annotations=$(cut -f 1,2,3 /Users/giacomo/Dropbox/11-Siboglinids
/05-PAPER/00-DATA/06-Metabolism/01-AA/expression/osedax_kallisto
_tpm_scaled_ok_ok.tsv | fgrep $line)
4   cat $annotations
5   kallisto_body=$(cut -f 2 <<< $annotations)
6   cat $kallisto_body
7   kallisto_roots=$(cut -f 3 <<< $annotations)
8   cat $kallisto_roots
```

```
9 echo $kallisto_body$\t'$kallisto_roots >> GlycineDegradation_kallisto_tpm_ALLscaled.tsv
10 done < GlycineDegradation_JUSTgeneIDs_osedax.txt
```

```
1 cut -f 1 GlycineDegradation_kallisto_tpm_ALLscaled.tsv | sed -e 's/ /\t/g' | cut -f 2,3 > GlycineDegradation_kallisto_tpm_ALLscaled_ONLY.tsv
2 cut -f 1 GlycineDegradation_geneIDs_osedax.txt > firstColumn
3 paste firstColumn GlycineDegradation_kallisto_tpm_ALLscaled_ONLY.tsv > FINAL_GlycineDegradation_osedax_kallisto_tpm_ALLscaled.tsv
```

```
1 library(tidyverse)
2 library(dplyr)
3 library(ggplot2)
4 library(data.table)
5 library(gplots)
6 library(pheatmap)
7 library(dendextend)
8 library(factoextra)
9 library(ComplexHeatmap)
10 library(RColorBrewer)
11 library(NbClust)
12 library(scales)
13
14 #Import data in matrix format
15 B <- read.delim("/Users/giacomo/Desktop/GlycineDegradation/FINAL_GlycineDegradation_osedax_kallisto_tpm_ALLscaled.tsv", row.names=1)
16
17 # To make 0 a different colour
18 # First create whatever gradient (e.g. RdBu)
```



```
19 heatmap_color <- colorRampPalette(brewer.pal(n = 7, name = "Red
  s"))(1000)
20 heatmap_color[1] <- rgb(1,1,1)
21 #column_labels = c("your","labels"),
22 #row_labels = c("your","labels"))
23
24 paletteLength <- 1000
25
26 pheatmap(B,
27           cluster_rows = FALSE,
28           cluster_cols = FALSE,
29           border_color = NA,
30           color = heatmap_color,
31           height = 5,
32           width = 25)
```

KEGG Metabolism

KofamKoala_universal_v1.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/metabolism
3 #$ -o /data/scratch/btx654/metabolism
4 #$ -j y
5 #$ -pe smp 12
6 #$ -l h_vmem=40G
7 #$ -l h_rt=72:0:0
8 #$ -l highmem
```

```
9
10 species=$1
11 NR_proteome="$species".fa
12 NR_proteome_path=/data/SBCS-MartinDuranLab/03-Giacomo/NR_proteomes/$NR_proteome
13 output_kofam="$species"_kofam_result.txt
14
15 echo "Working on "$species
16
17 module load anaconda3
18 conda activate kofam_env
19
20 mkdir $species
21 cd $species
22
23 exec_annotation \
24   --profile=/data/SBCS-MartinDuranLab/03-Giacomo/db/kofam/profiles/ \
25   --ko-list=/data/SBCS-MartinDuranLab/03-Giacomo/db/kofam/ko_list \
26   --cpu=12 \
27   --format=mapper \
28   --report-unannotated \
29   -o $output_kofam \
30   $NR_proteome_path
```

```
qsub KofamKoala_universal_v1.sh 0fra
```

select max 10000 sequences to submit to the online kofamkoala

```
1 awk -v RS='>' 'NR>1 { gsub("\n", ";", $0); sub(";$", "", $0); print ">"$0 }' 0fra.fa | head -n 10000 | tr ',' '\n' | sed "s/;/\n/g" > 0fra_first10000.fa
```

```

2  awk -v RS='>' 'NR>1 { gsub("\n", ";", $0); sub(";$", "", $0); pr
   int ">"$0 }' Ofra.fa | tail -n 8024 | tr ',' '\n' | sed "s/;/\n/
   g" > Ofra_last8024.fa

```

I got the results of kofamKOALA but they represent less protein compared to KAAS annotated ones: although I just noticed that we managed to annotate just 5522 proteins with kofamkoala whereas we annotated 8514 using KAAS so I will try to combine both the methods, kofamKOALA first then KAAS

```

1  awk '{print $2}' riftia_KAAS_custom_SBH.txt > only_K0numbers.txt
2  #add a line at the top of this file saying "K0_number"
3  nano only_K0numbers.txt
4  paste riftia_annotation_Dec2020_TrinoPanther.xls only_K0numbers.
   txt > riftia_annotation_Jan2021_TrinoPanther.xls
5
6  #awk '{print $1"\t"$2}' Ofra_kofam_result.txt > Ofra_kofam_resul
   t_ok.txt
7  cut -f 1 Ofra > Ofra_IDs_KAAS
8  grep "\sK" Ofra_kofam_result.txt | awk '{print $1"\t"$2}' | sed
   's/Ofra_//g' > Ofra_kofam_result_ok.txt
9  cut -f 1 Ofra_kofam_result_ok.txt > Ofra_IDs_kofam
10 #select the genes annotated with KAAS which don't have an annota
    tion with kofam
11 fgrep -vf Ofra_IDs_kofam Ofra | tail -n +2 > Ofra_KAAS_only
12 #combine the annotations using kofam as the primary one
13
14 cat Ofra_kofam_result_ok.txt Ofra_KAAS_only > Ofra_combined #932
   4 entries

```

combine.sh

```

1  #!/bin/bash

```

```
2
3 species=$1
4
5 grep "\sK" "$species"_kofam_result.txt | awk '{print $1"\t"$2}'
  | sed "s/$species//g" | sed "s/_//g" > "$species"_kofam_result_o
  k.txt
6 cut -f 1 "$species"_kofam_result_ok.txt > "$species"_IDs_kofam
7 #select the genes annotated with KAAS which don't have an annota
  tion with kofam
8 fgrep -vf "$species"_IDs_kofam ../$species | tail -n +2 > "$spe
  cies"_KAAS_only
9 #combine the annotations using kofam as the primary one
10 cat "$species"_kofam_result_ok.txt "$species"_KAAS_only > "$spec
  ies"_combined
```

```
cat Ofra_combined Rpac_combined Oalv_combined Lluy_combined Pech
_combined Ofus_combined Ctel_combined > K0list_OfraRpacOalvLluyP
echOfusCtel_combined.txt
```

reconstruct pathways

RESULTS
