

# New annotation Dec2020

## STEP 5

working directory:

/data/scratch/btx654/btx604-scratch/\$species/New\_annotation\_Dec2020/step5

in the working directory there must be:

- softmasked genome as "\$species"\_softmasked.fa
- mikado.loci.gff3
- Augustus output as "\$species".aug.out
- alignAssembly.config and annotCompare.config with the right paths specified

agat\_convert\_sp\_gxf2gxf.pl and fgrep

```
1 module load anaconda3
2 conda activate agat_env
3 agat_convert_sp_gxf2gxf.pl -g mikado.loci.gff3 -o mikado.loci.AG
  AT.gff3
4 #grep "ncRNA" mikado.loci.AGAT.gff3 | awk '{print $9}' | uniq >
  mikado.ncRNA.IDs
5 grep "ncRNA" mikado.loci.AGAT.gff3 | awk '{print $9}' | uniq | s
  ed "s=/\t/" | sed "s;/\t/" | awk '{print $2}' > mikado.ncRNA.I
  Ds
6 fgrep -v -w -f mikado.ncRNA.IDs mikado.loci.AGAT.gff3 > mikado.l
  oci.AGAT.NOncRNA.gff3
```

✓ oasisia

✓ osedax

✓ riftia

gffread\_universal\_5.1\_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=24:0:0
8
9  species=$1
10 species_softmasked="$species"_softmasked.fa
11 augustus_output="$species".aug.out
12 output_mikado_transcripts="$species"_mikado_transcripts_NoncRNA.
   fa
13 output_augustus_genes="$species".Augustus.genes.fa
14
15 echo "Working on "$species
16
17 module load anaconda3
18 source activate augustus
19
20 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
   ec2020/step5
21
22 gffread -C -x $output_mikado_transcripts -g $species_softmasked
   mikado.loci.AGAT.NoncRNA.gff3
23
24 gffread $augustus_output -V -w $output_augustus_genes -g $specie
   s_softmasked
```

✓ oasisia

✓ osedax

✓ riftia

seqclean\_universal\_5.2\_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=4G
7  #$ -l h_rt=24:0:0
8
9  species=$1
10 mikado_transcripts="$species"_mikado_transcripts_NoncRNA.fa
11
12 echo "Working on "$species
13
14 module load anaconda3
15 source activate pasa
16
17 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
  ec2020/step5
18
19 /data/home/btx654/.conda/envs/pasa/opt/pasa-2.4.1/bin/seqclean
  $mikado_transcripts
```

✓ oasisia

✓ osedax

✓ riftia

#### pasa\_universal\_5.3\_v1.highmem.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 24
6  #$ -l h_vmem=25G
```

```
7  #$ -l h_rt=240:0:0
8  #$ -l highmem
9
10 species=$1
11 species_softmasked="$species"_softmasked.fa
12 mikado_transcripts="$species"_mikado_transcripts_NoncRNA.fa
13 mikado_transcripts_clean="$species"_mikado_transcripts_NoncRNA.f
   a.clean
14
15 echo "Working on "$species
16
17 module load anaconda3
18 source activate pasa
19 module load samtools
20
21 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
   ec2020/step5
22
23 /data/home/btx654/.conda/envs/pasa/opt/pasa-2.4.1/Launch_PASA_pi
   peline.pl -c alignAssembly.config -C -R -g $species_softmasked -
   t $mikado_transcripts_clean -T -u $mikado_transcripts --ALIGNERS
   blat --CPU 24
```

✓ oasisia

✓ osedax

✓ riftia

#### pasa\_universal\_5.4\_v1.sh

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

```
6  # $ -l h_vmem=10G
7  # $ -l h_rt=48:0:0
8
9  species=$1
10 species_softmasked="$species"_softmasked.fa
11 augustus_output="$species".aug.out
12 augustus_gtf="$species".aug.gtf
13
14 echo "Working on "$species
15
16 module load anaconda3
17 source activate augustus
18
19 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_De2020/step5
20
21 gffread $augustus_output -E -T -o $augustus_gtf
22
23 conda deactivate
24 source activate pasa
25 module load samtools/1.9
26
27 /data/home/btx654/.conda/envs/pasa/opt/pasa-2.4.1/scripts/Load_Current_Gene_Annotations.dbi -c alignAssembly.config -g $species_softmasked -P $augustus_gtf
```

- ☒ oasisia
- ☒ osedax
- ☒ riftia

pasa\_universal\_5.5\_v1.highmem.sh

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

```
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 30
6  #$ -l h_vmem=3G
7  #$ -l h_rt=48:0:0
8  #$ -l highmem
9
10 species=$1
11 species_softmasked="$species"_softmasked.fa
12 mikado_transcripts_clean="$species"_mikado_transcripts_NoncRNA.f
   a.clean
13
14 echo "Working on "$species
15
16 module load anaconda3
17 source activate pasa
18 module load samtools
19
20 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
   ec2020/step5
21
22 /data/home/btx654/.conda/envs/pasa/opt/pasa-2.4.1/Launch_PASA_pi
   peline.pl -c annotCompare.config -A -g $species_softmasked -t $m
   ikado_transcripts_clean --CPU 30
```

✓ oasisia

✓ osedax

✓ riftia

change name pasa gff3 output

```
1 mv sqlite_db.gene_structures_post_PASA_updates.16689.gff3 oasisi
   a_pasa_FirstStep.gff3
```

```
2 mv sqlite_db.gene_structures_post_PASA_updates.20574.gff3 riftia
   _pasa_FirstStep.gff3
3 mv sqlite_db.gene_structures_post_PASA_updates.47524.gff3 osedax
   _pasa_FirstStep.gff3
```

✓ oasisia

✓ osedax

✓ riftia

#### pasa\_universal\_5.6\_v1.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/
3 #$ -o /data/scratch/btx654/
4 #$ -j y
5 #$ -pe smp 2
6 #$ -l h_vmem=10G
7 #$ -l h_rt=48:0:0
8
9 species=$1
10 species_softmasked="$species"_softmasked.fa
11 pasa_first_step_gff3="$species"_pasa_FirstStep.gff3
12 pasa_first_step_gtf="$species"_pasa_FirstStep.gtf
13
14
15 module load anaconda3
16 source activate augustus
17
18
19 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
   ec2020/step5
20
21 gffread $pasa_first_step_gff3 -E -T -o $pasa_first_step_gtf
22
```

```
23 conda deactivate
24 source activate pasa
25 module load samtools/1.9
26
27 /data/home/btx654/.conda/envs/pasa/opt/pasa-2.4.1/scripts/Load_C
  urrent_Gene_Annotations.dbi -c alignAssembly.config -g $species_
  softmasked -P $pasa_first_step_gtf
```

✓ oasisia

✓ osedax

✓ riftia

Repeating "pasa\_universal\_5.5\_v1.highmem.sh"

✓ oasisia

✓ osedax

✓ riftia

change name pasa gff3 and bed outputs

```
1 mv sqlite_db.gene_structures_post_PASA_updates.45674.gff3 riftia
  _pasa_SecondStep.gff3
2 mv sqlite_db.gene_structures_post_PASA_updates.45674.bed riftia_
  pasa_SecondStep.bed
3 cp riftia_pasa_SecondStep.* /data/SBCS-MartinDuranLab/03-Giacomo
  /data/riftia/annotation/New_annotation_Dec2020/step5/
4 mv sqlite_db.gene_structures_post_PASA_updates.37064.gff3 oasisi
  a_pasa_SecondStep.gff3
5 mv sqlite_db.gene_structures_post_PASA_updates.37064.bed oasisia
  _pasa_SecondStep.bed
6 cp oasisia_pasa_SecondStep.* /data/SBCS-MartinDuranLab/03-Giacom
  o/data/oasisia/annotation/New_annotation_Dec2020/step5/
7 mv sqlite_db.gene_structures_post_PASA_updates.11480.gff3 osedax
  _pasa_SecondStep.gff3
8 mv sqlite_db.gene_structures_post_PASA_updates.11480.bed osedax_
  pasa_SecondStep.bed
9 cp osedax_pasa_SecondStep.* /data/SBCS-MartinDuranLab/03-Giacomo
```



```
/data/osedax/annotation/New_annotation_Dec2020/step5/
```

- ✓ oasisia
- ✓ osedax
- ✓ riftia

## STEP 6

create a folder which can be deleted after this step6:

```
/data/SBCS-MartinDuranLab/03-Giacomo/data/$species/annotation  
/repeatmasker_delete_after_step6_Dec2020
```

containing the output of repeatmasker named as:

- "\$species"\_repeatmasker.fa.out

filtering\_universal\_6.1\_v1.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 2
6  #$ -l h_vmem=3G
7  #$ -l h_rt=24:0:0
8
9  species=$1
10 species_softmasked="$species"_softmasked.fa
11 species_softmasked_path=/data/scratch/btx654/btx604-scratch/$spe
    cies/New_annotation_Dec2020/step5/$species_softmasked
12 pasa_second_step_gff3="$species"_pasa_SecondStep.gff3
13 pasa_second_step_path=/data/scratch/btx654/btx604-scratch/$speci
    es/New_annotation_Dec2020/step5/$pasa_second_step_gff3
14 pasa_second_step_AGAT="$species"_pasa_SecondStep.AGAT.gff3
15 output_gff3="$species".AGAT.noSTOP.gff3
```

```
16
17 echo "Working on "$species
18
19 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
ec2020/
20
21 mkdir -p step6
22 cd step6
23
24 cp $species_softmasked_path ./
25 cp $pasa_second_step_path ./
26
27 module load anaconda3
28 conda activate agat_env
29
30 agat_convert_sp_gxf2gxf.pl -g $pasa_second_step_gff3 -o $pasa_se
cond_step_AGAT
31
32 conda deactivate
33 source activate augustus
34
35 gffread -E $pasa_second_step_AGAT -g $species_softmasked -V -H -
o $output_gff3
```

- ✓ oasisia
- ✓ osedax
- ✓ riftia

```
1 grep -c "^Warning: In-frame STOP found for" filtering_universal_
6.1_v1.sh.o1397518 #291 riftia
2 grep -c "^Warning: In-frame STOP found for" filtering_universal_
6.1_v1.sh.o1397519 #451 oasisia
3 grep -c "^Warning: In-frame STOP found for" filtering_universal_
```

```
6.1_v1.sh.o1399077 #232 osedax
```

### filtering\_universal\_6.2\_v1.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 2
6  #$ -l h_vmem=5G
7  #$ -l h_rt=24:0:0
8
9  species=$1
10 pasa_gff3="$species".AGAT.noSTOP.gff3
11 pasa_gtf="$species".AGAT.noSTOP.gtf
12 repeatmasker="$species"_repeatmasker.fa.out
13 repeatmasker_path=/data/SBCS-MartinDuranLab/03-Giacomo/data/$species/annotation/repeatmasker_delete_after_step6_Dec2020/$repeatmasker
14 output_filt_gtf="$species".AGAT.noSTOP.filt.gtf
15 output_filt_gff3="$species".AGAT.noSTOP.filt.AGAT.gff3
16 species_softmasked="$species"_softmasked.fa
17
18 echo "Working on "$species
19
20 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_Dec2020/step6/
21
22 module load anaconda3
23 source activate augustus
24
```

```
25 gffread -E $pasa_gff3 -T -o $pasa_gtf
26
27 conda deactivate
28 module unload anaconda3
29
30 module load python/2.7.15
31 module load bedtools
32
33 cp $repeatmasker_path ./
34 sed -e 's/.arrow.arrow.pilon.pilon//' $repeatmasker > repeatmask
er_sed.fa.out
35
36 python2 /data/SBCS-MartinDuranLab/03-Giacomo/src/various/rep2be
d.py repeatmasker_sed.fa.out > RepeatMasker.bed
37 python2 /data/SBCS-MartinDuranLab/03-Giacomo/src/various/filt-re
p-gtf.py $pasa_gtf RepeatMasker.bed
38
39 module unload python/2.7.15
40 module unload bedtools
41 module load anaconda3
42 conda activate agat_env
43
44 agat_convert_sp_gxf2gxf.pl -g $output_filt_gtf -o $output_filt_g
ff3
45 agat_sq_stat_basic.pl -i $output_filt_gff3 -g $species_softmaske
d
```

✓ oasisia

✓ osedax

✓ riftia

## Results:

- Riftia

238207 exons...

initial number of genes: 43443

number of genes after filtering 37455

Type (3rd column)	Number	Size total (kb)	Size mean (bp)	% of the genome
-------------------	--------	-----------------	----------------	-----------------

Results are rounded to two decimal places

cds	226079	55718.80	246.46	10.07
-----	--------	----------	--------	-------

exon	226211	55876.82	247.01	10.09
------	--------	----------	--------	-------

five_prime_utr	1531	91.68	59.89	0.02
----------------	------	-------	-------	------

gene	37455	310501.42	8289.99	56.09
------	-------	-----------	---------	-------

three_prime_utr	236	66.33	281.06	0.0
-----------------	-----	-------	--------	-----

1

transcript	38594	342050.72	8862.80	6
------------	-------	-----------	---------	---

1.79

Total	530106	764305.77	1441.80	138.0
-------	--------	-----------	---------	-------

7

#### • Oasisia

337078 exons...

initial number of genes: 62270

number of genes after filtering 37929

Type (3rd column)	Number	Size total (kb)	Size mean (bp)	% of the genome
-------------------	--------	-----------------	----------------	-----------------

Results are rounded to two decimal places

cds	279272	62067.95	222.25	7.68
-----	--------	----------	--------	------

exon	279538	62212.03	222.55	7.70
------	--------	----------	--------	------

five_prime_utr	1993	102.59	51.47	0.0
----------------	------	--------	-------	-----

1

gene	37929	419731.45	11066.24	51.95
------	-------	-----------	----------	-------

three_prime_utr	125	41.49	331.94	0.0
-----------------	-----	-------	--------	-----

1

11	transcript	39850	482789.42	12115.17	59.76
12	Total	638707	1026944.94	1607.85	127.11

• Osedax

1	185496 exons...				
2	initial number of genes: 21969				
3	number of genes after filtering 18176				
4					
5	Type (3rd column)	Number	Size total (kb)	Size mean (bp)	% of the genome
	/!\Results are rounding to two decimal places				
6	cds	174271	29205.60	167.59	10.26
7	exon	174426	29269.47	167.80	10.28
8	five_prime_utr	1278	47.57	37.22	0.02
9	gene	18176	173417.61	9541.02	60.92
10	three_prime_utr	32	16.30	509.44	0.01
11	transcript	18808	183673.13	9765.69	64.52
12	Total	386991	415629.70	1074.00	146.00

filtering\_universal\_6.3\_v1.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=5G
7  #$ -l h_rt=24:0:0
8
9  species=$1

```

```
10 pasa_gtf="$species".AGAT.noSTOP.filt.gtf
11 species_softmasked="$species"_softmasked.fa
12 pasa_prot_fasta="$species".AGAT.noSTOP.filt.prot.fasta
13 final_pasa_gtf="$species".AGAT.noSTOP.filt.noTE.gtf
14 final_pasa_prot_fasta="$species".AGAT.noSTOP.filt.noTE.prot.fasta
15 final_pasa_gff3="$species".AGAT.noSTOP.filt.noTE.AGAT.gff3
16
17 echo "Working on "$species
18
19 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_Dec2020/step6/
20
21 module load anaconda3
22 source activate augustus
23
24 gffread -E $pasa_gtf -S -g $species_softmasked -y $pasa_prot_fasta
25
26 conda deactivate
27 module unload anaconda3
28
29 cp /data/SBCS-MartinDuranLab/03-Giacomo/src/RepeatMasker/Libraries/RepeatPeps.lib ./
30 module load diamond/0.9.22
31 diamond makedb --in RepeatPeps.lib -d RepeatPeps
32 diamond blastp -d RepeatPeps -q $pasa_prot_fasta -o pasa.vs.RepeatPeps.1e5.blastp -f 6 qseqid bitscore evalue stitle -k 25 -e 1e-5 -p 8
33
34 cat pasa.vs.RepeatPeps.1e5.blastp | cut -f 1 | sort | uniq > TESIDs.txt
```

```
35 fgrep -w -v -f TEsIDs.txt $pasa_gtf > $final_pasa_gtf
36
37 #check if we still have TEs in our proteins
38 module load anaconda3
39 source activate augustus
40
41 gffread -E $final_pasa_gtf -S -g $species_softmasked -y $final_p
asa_prot_fasta
42
43 conda deactivate
44 module unload anaconda3
45
46 module load diamond/0.9.22
47 diamond blastp -d RepeatPeps -q $final_pasa_prot_fasta -o pasa_n
oTE.vs.RepeatPeps.1e5.blastp -f 6 qseqid bitscore evalue stitle
-k 25 -e 1e-5 -p 8
48 #check end
49
50 module load anaconda3
51 conda activate agat_env
52
53 agat_convert_sp_gxf2gxf.pl -g $final_pasa_gtf -o $final_pasa_gff
3
54 agat_sq_stat_basic.pl -i $final_pasa_gff3 -g $species_softmasked
```

- ✓ oasisia
- ✓ osedax
- ✓ riftia

#### Results:

- Riftia

```
1 415 queries aligned.
```

```
2
```



Type (3rd column)	Number	Size total (kb)	Si
ze mean (bp)	% of the genome	Results are roundi	ng to two decimal places
cds	224328	55173.67	245.95 9.97
exon	224460	55330.47	246.50 10.00
five_prime_utr	1520	90.92	59.81 0.02
gene	37043	307881.72	8311.47 55.62
three_prime_utr	233	65.88	282.76 0.0
transcript	38179	339383.55	8889.27 6
Total	525763	757926.20	1441.57 136.9

- Oasisia

1	2072 queries aligned.				
2					
3	Type (3rd column)	Number	Size total (kb)	Size mean (bp) % of the genome	
	/!\Results are rounding to two decimal places				
4	cds	270655	58746.82	217.05	7.27
5	exon	270916	58885.90	217.36	7.29
6	five_prime_utr	1970	100.30	50.91	0.0
7	gene	35869	405225.34	11297.37	50.16
8	three_prime_utr	116	38.78	334.31	0.0
9	transcript	37777	467927.45	12386.57	57.92
10	Total	617303	990924.59	1605.25	122.6

- Osedax

151 queries aligned.
----------------------

3	Type (3rd column)	Number	Size total (kb)	Si
	ze mean (bp)	% of the genome	/\Results are roundi	
	ng to two decimal places			
4	cds	173041	28963.63	167.38
5	exon	173196	29026.81	167.60
6	five_prime_utr	1276	47.55	37.27
7	gene	18025	172119.65	9548.94
8	three_prime_utr	30	15.63	520.97
9	transcript	18657	182375.17	9775.16
				6
				4.07
10	Total	384225	412548.44	1073.72
				144.9
				2

#### rename\_and\_longest\_isoform.sh

```

1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=5G
7  #$ -l h_rt=2:0:0
8
9  species=$1
10 prefix=$2
11 final_pasa_gff3="$species".AGAT.noSTOP.filt.noTE.AGAT.gff3
12 final_pasa_gtf="$species".AGAT.noSTOP.filt.noTE.gtf
13 output_annotation="$species"_annotation_v101220.gff3
14 loci_merged="$species"_lociMerged.gff
15 longest_isoform="$species"_lociMerged_longestIsoform.gff
16
17 echo "Working on "$species" using prefix: "$prefix

```

```
18
19 mkdir /data/SBCS-MartinDuranLab/03-Giacomo/data/$species/annotation/New_annotation_Dec2020/step6
20 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_Dec2020/step6/
21
22 cp $final_pasa_gff3 /data/SBCS-MartinDuranLab/03-Giacomo/data/$species/annotation/New_annotation_Dec2020/step6/
23 cp $final_pasa_gtf /data/SBCS-MartinDuranLab/03-Giacomo/data/$species/annotation/New_annotation_Dec2020/step6/
24
25 cd /data/SBCS-MartinDuranLab/03-Giacomo/data/$species/annotation/New_annotation_Dec2020/step6/
26
27 module load anaconda3
28 conda activate agat_env
29
30 agat_sp_manage_IDs.pl -f $final_pasa_gff3 --ensembl --prefix $prefix --type_dependent --tair -o $output_annotation
31
32 agat_convert_sp_gxf2gxf.pl --gff $output_annotation --merge_loci -o $loci_merged
33 agat_sp_keep_longest_isoform.pl --gff $loci_merged -o $longest_isoform
```

- ✓ ~~oasisia~~ OALV
- ✓ ~~osedax~~ OFRA
- ✓ ~~riftia~~ RPAG

## STEP 7

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=$1
11 annotation_gtf="$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/btx604-scratch/$species/New_annotation_De2020/step6/
19
20 module load anaconda3
21 source activate augustus
22
23 gffread -E $annotation_gtf -g $species_softmasked -y $annotation_fa
24
25 conda deactivate
26 source activate busco_env
27 #export BUSCO_CONFIG_FILE="/data/home/btx654/.conda/envs/busco_env/busco/config/myconfig.ini"
28 #export AUGUSTUS_CONFIG_PATH=/data/SBCS-MartinDuranLab/02-Chema/src/Augustus/config/
29
```

```
30 busco -i $annotation_fa -m proteins -o $output_busco -c 4 -l met
    azoa_odb10
31
32 cd $output_busco
33 cd run_*
34 mkdir /data/SBCS-MartinDuranLab/03-Giacomo/data/$species/annotat
    ion/New_annotation_Dec2020/step7
35 cp full_table.tsv /data/SBCS-MartinDuranLab/03-Giacomo/data/$spe
    cies/annotation/New_annotation_Dec2020/step7
36 cp missing_busco_list.tsv /data/SBCS-MartinDuranLab/03-Giacomo/d
    ata/$species/annotation/New_annotation_Dec2020/step7
37 cp short_summary.txt /data/SBCS-MartinDuranLab/03-Giacomo/data
    /$species/annotation/New_annotation_Dec2020/step7
☒ oasisia
☒ esedax
☒ riftia
```

## STEP 8

gffread\_universal\_8.1\_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=24:0:0
8  #$ -l highmem
9
10 species=$1
11 final_annotation="$species"_annotation_*.gff3
```

```
12 final_annotation_path=/data/SBCS-MartinDuranLab/03-Giacomo/data
   /$species/annotation/New_annotation_Dec2020/step6/$final_annotat
   ion
13 species_softmasked="$species"_softmasked.fa
14 species_softmasked_path=/data/scratch/btx654/btx604-scratch/$spe
   cies/New_annotation_Dec2020/step6/$species_softmasked
15 output_mRNA="$species"_mRNA.fa
16 output_CDS="$species"_CDS.fa
17 output_proteins="$species"_proteins.fa
18
19 echo "Working on "$species
20
21 module load anaconda3
22 source activate augustus
23
24 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
   ec2020/
25
26 mkdir -p step8
27 cd step8
28
29 cp $species_softmasked_path ./
30 cp $final_annotation_path ./
31
32 gffread -w $output_mRNA -g $species_softmasked $final_annotation
33 gffread -x $output_CDS -g $species_softmasked $final_annotation
34 gffread -y $output_proteins -g $species_softmasked $final_annota
   tion
```

- ✓ oasisia
- ✓ osedax
- ✓ riftia

## obtain\_trinotate\_inputs\_8.1\_v1.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=20G
7  #$ -l h_rt=72:0:0
8  #$ -l highmem
9
10 species=$1
11 fasta_mRNA="$species"_mRNA.fa
12 fasta_CDS="$species"_CDS.fa
13 fasta_proteins="$species"_proteins.fa
14 gene_trans_map="$species".gene_trans_map
15
16 echo "Working on "$species
17
18 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_De2020/step8/
19
20 grep ">" $fasta_mRNA | awk '{ print $2"(")}' > positions.txt
21 sed -i 's/CDS=/:/' positions.txt
22
23 awk '/^>/ {if (seqlen){print "len:"seqlen}; print ;seqlen=0;next; } { seqlen += length($0)}END{print "len:"seqlen}' $fasta_proteins | grep -v ">" > lengths.txt
24
25 grep ">" $fasta_proteins > names.txt
26 sed 's/>/' names.txt > names_clean.txt
27
```

```
28 sed '/^>/s/ .*//' $fasta_mRNA > input_trinotate_mRNA.fa
29
30 cp $fasta_proteins input_trinotate_proteins.fa
31 INDEX=1
32 while read -r line
33 do
34     original_name=$(echo $line)
35     lenght=$(head -$INDEX lenghts.txt | tail -1)
36     name_clean=$(head -$INDEX names_clean.txt | tail -1)
37     position=$(head -$INDEX positions.txt | tail -1)
38
39     sed -i "s/$original_name/$original_name $lenght $name_clean$position/" input_trinotate_proteins.fa
40
41     INDEX=$((INDEX+1))
42 done < names.txt
43
44
45 grep ">" input_trinotate_proteins.fa > full_names.txt
46 sed -i 's/>/' full_names.txt
47
48 INDEX=1
49 while read -r line
50 do
51     full_name=$(echo $line)
52     name_clean=$(head -$INDEX names_clean.txt | tail -1)
53
54     echo -e $full_name'\t'$name_clean >> $gene_trans_map
55
56     INDEX=$((INDEX+1))
```



```
57 done < full_names.txt
```

- ✓ oasisia
- ✓ osedax
- ✓ riftia

### 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=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 -evaluate 1e-3 > blastp.outfmt6
```

### BLASTx\_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=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  blastx -query input_trinotate_mRNA.fa -db /data/SBCS-MartinDuran
    Lab/03-Giacomo/db/trinotate/uniprot_sprot.pep -num_threads 8 -ma
    x_target_seqs 1 -outfmt 6 -evaluate 1e-3 > blastx.outfmt6
```

### HMMER\_universal.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/
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 species=$1
12
13 echo "Working on "$species
14
15 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
  ec2020/step8/
16
17 module load anaconda3
18 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
  3/trinotate_env
19
20 hmmscan --cpu 12 --domtblout PFAM.out /data/SBCS-MartinDuranLab/
  03-Giacomo/db/trinotate/Pfam-A.hmm input_trinotate_proteins.fa >
  pfam.log
```

- i got "Segmentation fault" for oasisia and osedax so I will send this job with more ram (40G) for them

---

This section is intended to fix the "Segmentation Fault" error of the previous script. Basically we split the input\_trinotate\_proteins.fa in 100 parts in order to not overload HMMER

obtain\_rename\_chunks\_universal\_v1.sh

```
1 #!/bin/bash
2 #$ -cwd
3 #$ -pe smp 1
4 #$ -l h_vmem=1G
5 #$ -j y
6 #$ -l h_rt=01:00:00
7 species=$1
8 trinotate_folder=input_trinotate_proteins_chunks
9 trinotate_chunk=input_trinotate_proteins.fa_chunk_
```

```
10
11 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
    ec2020/step8/
12
13 mkdir $trinotate_folder
14 module load exonerate/2.4.0
15 fastasplit -f input_trinotate_proteins.fa -o ./input_trinotate_p
    roteins_chunks/ -c 100
16 cd $trinotate_folder
17
18 readlink -f "$trinotate_chunk"* > list
19 for i in $(seq 1 100);
20 do
21     original_file=$(head -"$i" list | tail -1)
22     renamed_file="$trinotate_chunk"$i
23     mv -- "$original_file" "$renamed_file"
24 done
```

Run HMMER on the single chunks:

HMMER\_universal\_chunks.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=40G
7 #$ -l h_rt=36:0:0
8 #$ -t 1-100
9 #$ -l highmem
10
11 species=$1
```

```
12 target_chunk=input_trinotate_proteins.fa_chunk_"${SGE_TASK_ID}"
13 PFAM_out=PFAM_"${SGE_TASK_ID}".out
14 pfam_log=pfam_"${SGE_TASK_ID}".log
15
16 echo "Working on "$species
17
18 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
  ec2020/step8/input_trinotate_proteins_chunks
19
20 module load anaconda3
21 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
  3/trinotate_env
22
23 hmmscan --cpu 12 --domtblout $PFAM_out /data/SBCS-MartinDuranLab
  /03-Giacomo/db/trinotate/Pfam-A.hmm $target_chunk > $pfam_log
```

Some jobs will crash because of the same error so we can divide in 100 subchunks the chunks that failed the previous jobs. With these steps we will remove those problematic subchunks from the analysis. Every subchunk removed in this way will be just the 0.01% of the total so it won't impact too much the consistency of our analyses

obtain\_rename\_chunks\_universal\_fix\_v1.sh

```
1 #!/bin/bash
2 #$ -cwd
3 #$ -pe smp 1
4 #$ -l h_vmem=1G
5 #$ -j y
6 #$ -l h_rt=01:00:00
7
8 species=$1
9 chunk=$2
```

```
10 trinotate_folder_chunk=input_trinotate_proteins_chunk_"$chunk"
11 chunk_fix=input_trinotate_proteins.fa_chunk_"$chunk"
12 trinotate_chunk=input_trinotate_proteins.fa_chunk_
13
14 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
ec2020/step8/input_trinotate_proteins_chunks/
15
16 mkdir $trinotate_folder_chunk
17 module load exonerate/2.4.0
18 fastasplit -f $chunk_fix -o ./$trinotate_folder_chunk/ -c 100
19 cd $trinotate_folder_chunk
20
21 readlink -f "$trinotate_chunk"* > list
22 for i in $(seq 1 100);
23 do
24     original_file=$(head -$i list | tail -1)
25     renamed_file="$trinotate_chunk"$i
26     mv -- "$original_file" "$renamed_file"
27 done
```

Run HMMER on the subchunks

HMMER\_universal\_chunks\_fix\_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=1G
7 #$ -l h_rt=36:0:0
8 #$ -t 1-100
9 #$ -l highmem
```

```
10
11 species=$1
12 chunk=$2
13 trinotate_folder_chunk=input_trinotate_proteins_chunk_"$chunk"
14 target_chunk=input_trinotate_proteins.fa_chunk_"${SGE_TASK_ID}"
15 PFAM_out=PFAM_"${SGE_TASK_ID}".out
16 pfam_log=pfam_"${SGE_TASK_ID}".log
17
18 echo "Working on "$species" on chunk: "$chunk
19
20 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
ec2020/step8/input_trinotate_proteins_chunks/$trinotate_folder_c
hunk
21
22 module load anaconda3
23 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
3/trinotate_env
24
25 hmmscan --cpu 12 --domtblout $PFAM_out /data/SBCS-MartinDuranLab
/03-Giacomo/db/trinotate/Pfam-A.hmm $target_chunk > $pfam_log
```

to be removed:

- oasisia 38 37 rm PFAM\_37.out
- oasisia 46 25 rm PFAM\_25.out
- oasisia 75 32 rm PFAM\_32.out
- osedax 75 81 rm PFAM\_81.out
- osedax 52 70 rm PFAM\_70.out

Now we need to put together the PFAM.out file from the subchunks to basically rebuild the PFAM.out file of the chunk that failed before

merge\_chunks\_universal\_fix\_v1.sh

```
1 #!/bin/bash
2 # $ -cwd
3 # $ -pe smp 1
```

```
4  #$ -l h_vmem=1G
5  #$ -j y
6  #$ -l h_rt=01:00:00
7
8  species=$1
9  chunk=$2
10 trinotate_folder_chunk=input_trinotate_proteins_chunk_"$chunk"
11 PFAM_chunk_out=PFAM_"$chunk".out
12
13 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
ec2020/step8/input_trinotate_proteins_chunks/$trinotate_folder_c
hunk
14
15 for i in $(seq 1 100);
16 do
17     PFAM_out=PFAM_"$i".out
18     PFAM_out_ok=PFAM_"$i"_ok.out
19     head -n -10 $PFAM_out | tail -n +4 > $PFAM_out_ok
20 done
21 cat PFAM_*_ok.out > ../$PFAM_chunk_out
```

And finally we need to put together the PFAM.out file of all the chunks to reconstruct the PFAM.out file that we would have obtained if we didn't encounter the Segmentation fault error

merge\_chunks\_universal\_v1.sh

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -pe smp 1
4  #$ -l h_vmem=1G
5  #$ -j y
6  #$ -l h_rt=01:00:00
```



```
7
8 species=$1
9
10 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
    ec2020/step8/input_trinotate_proteins_chunks
11
12 for i in $(seq 1 100);
13 do
14     PFAM_out=PFAM_"$i".out
15     PFAM_out_ok=PFAM_"$i"_ok.out
16     head -n -10 $PFAM_out | tail -n +4 > $PFAM_out_ok
17 done
18 cat PFAM*_ok.out > ../PFAM.out
```

End of fix

---

#### Signalp\_universal.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/
3 #$ -o /data/scratch/btx654/
4 #$ -j y
5 #$ -pe smp 1
6 #$ -l h_vmem=8G
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 perl
16
17 /data/SBCS-MartinDuranLab/03-Giacomo/src/signalp-4.1/signalp -f
    short -n signalp.out input_trinotate_proteins.fa
```

#### wrapper\_trinotate\_8.2\_v1.sh

```
1 #!/bin/bash
2 #$ -wd /data/home/btx654/scripts/annotation/New_annotation_Dec20
    20/step8/
3 #$ -j y
4 #$ -pe smp 1
5 #$ -l h_vmem=1G
6 #$ -l h_rt=1:0:0
7
8 species=$1
9
10 echo "Working on "$species
11
12 qsub BLASTp_universal.sh $species
13 qsub BLASTx_universal.sh $species
14 qsub HMMER_universal.sh $species
15 qsub Signalp_universal.sh $species
```

- ✓ oasisia
- ✓ osedax
- ✓ riftia

#### trinotate\_SQLite\_database\_v1.sh

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

```
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 1
6  #$ -l h_vmem=8G
7  #$ -l h_rt=2:0:0
8
9  species=$1
10 gene_trans_map="$species".gene_trans_map
11 sqlite_db="$species".sqlite
12
13 echo "Working on "$species
14
15 module load anaconda3
16 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
17 3/trinotate_env
18
19 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
20 ec2020/step8/
21
22 Build_Trinotate_Boilerplate_SQLite_db.pl $species
23
24 Trinotate $sqlite_db init --gene_trans_map $gene_trans_map --tra
25 nscript_fasta input_trinotate_mRNA.fa --transdecoder_pep input_t
26 rinotate_proteins.fa
27
28 #LOAD annotations, below are examples from Trinotate manual, cha
29 nge accordingly to your species ouput from blast+:
30
31 Trinotate $sqlite_db LOAD_swissprot_blastp blastp.outfmt6
32 Trinotate $sqlite_db LOAD_swissprot_blastx blastx.outfmt6
33 Trinotate $sqlite_db LOAD_signalp signalp.out
34 Trinotate $sqlite_db LOAD_pfam PFAM.out
35
```

```
29 Trinotate $sqlite_db report --incl_pep --incl_trans > annotation
_report.xls
```

- ✓ oasisia
- ✓ osedax
- ✓ riftia

## STEP 9

### pantherScore\_universal\_v1.sh

```
1  #!/bin/bash
2  #$ -wd /data/scratch/btx654/
3  #$ -o /data/scratch/btx654/
4  #$ -j y
5  #$ -pe smp 30
6  #$ -l h_vmem=5G
7  #$ -l h_rt=240:00:0
8  #$ -l highmem
9
10 species=$1
11 gffread_proteins="$species"_proteins.fa #generated by gffread_un
    iversal_8.1_v1.sh
12 gffread_proteins_path=/data/scratch/btx654/btx604-scratch/$speci
    es/New_annotation_Dec2020/step8/$gffread_proteins
13 panther_output="$species"_Panther
14
15 echo "Working on "$species
16
17 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
    ec2020/
18
```

```
19 mkdir -p step9
20 cd step9
21
22 cp $gffread_proteins_path ./
23
24 module load perl
25 module load hmmer/
26
27 export PERL5LIB=/data/SBCS-MartinDuranLab/03-Giacomo/src/hmmscoring/lib/
28
29 perl /data/SBCS-MartinDuranLab/03-Giacomo/src/hmmscoring/pantherScore2.2.pl -l /data/SBCS-MartinDuranLab/03-Giacomo/src/hmmscoring/PANTHER15.0/ -D B -n -o $panther_output -i $gffread_proteins -c 30 -V -s
30
31 mkdir /data/SBCS-MartinDuranLab/03-Giacomo/data/$species/annotation/New_annotation_Dec2020/step9
32 cp "$species"_Panthe* /data/SBCS-MartinDuranLab/03-Giacomo/data/$species/annotation/New_annotation_Dec2020/step9
```

✓ oasisia

✓ osedax

✓ riftia

#### combining Panther and Trinotate:

```
1 cd /data/scratch/btx654/btx604-scratch/riftia/annotation_step9
2 cp /data/SBCS-MartinDuranLab/03-Giacomo/data/riftia/annotation/step8/riftia_annotation_report.xls ./
3 sort oasisia_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 -v -f IDs_panther IDs_all > IDs_absentPanther ### There are oasisia:8632 osedax:4052 riftia:8182 genes without Panther ann
```

```

otation
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          sprotop_BLASTX_hit
   RNAMMER            prot_id            prot_coords          sprotop_BL
   ASTP_hit           Pfam            SignalP            TmHMM            eggnog
   Kegg              gene_ontology_BLASTX          gene_ontology_BLAST
   P                gene_ontology_Pfam          transcript          peptide
13 paste annotation_report.xls Panther_sorted_allgenes_noduplicates
   > oasisia_annotation_Dec2020_TrinoPanther.xls

```

- ✓ oasisia
- ✓ osedax
- ✓ riftia

## KofamKOALA

```

1  module load anaconda3
2  conda create --prefix /data/SBCS-MartinDuranLab/03-Giacomo/src/a
   naconda3/KofamKOALA_env
3  conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
   3/KofamKOALA_env
4  conda install -c bioconda kofamscan
5
6  cd /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/KofamKOALA
   _env
7  wget -r ftp://ftp.genome.jp/pub/db/kofam/profiles.tar.gz
8  cd ftp.genome.jp/pub/db/kofam
9  tar -zxvf profiles.tar.gz

```

```
10 wget ftp://ftp.genome.jp/pub/db/kofam/ko_list.gz
11 gzip -d ko_list.gz
```

### KofamKoala\_universal\_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=40G
7  #$ -l h_rt=72:0:0
8  #$ -l highmem
9
10 species=$1
11 gffread_proteins="$species"_proteins.fa #generated by gffread_un
    iversal_8.1_v1 .sh
12 gffread_proteins_path=/data/scratch/btx654/btx604-scratch/$speci
    es/New_annotation_Dec2020/step8/$gffread_proteins
13 output_kofam="$species"_kofam_result.txt
14
15 echo "Working on "$species
16
17 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
    ec2020/step9
18 mkdir -p kofamKoala
19 cd kofamKoala
20
21 cp $gffread_proteins_path ./
22 awk '/^>/ {printf("\n%s\n",$0);next; } { printf("%s",$0);} END
    {printf("\n");}' < $gffread_proteins | tail -n +2 > single_fasta
    _proteins.fa
```

```
23
24 module load anaconda3
25 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
3/KofamKOALA_env
26
27 exec_annotation \
28 --profile=/data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/Ko
famKOALA_env/ftp.genome.jp/pub/db/kofam/profiles/ \
29 --ko-list=/data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/Ko
famKOALA_env/ftp.genome.jp/pub/db/kofam/ko_list \
30 --cpu=12 \
31 --format=mapper \
32 --report-unannotated \
33 -o $output_kofam \
34 single_fasta_proteins.fa
```

✓ oasisia

✓ osedax

✓ riftia

I think I have to divide in chunks

obtain\_rename\_chunks\_universal\_v1.sh

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/
3 #$ -o /data/scratch/btx654/
4 #$ -pe smp 1
5 #$ -l h_vmem=1G
6 #$ -j y
7 #$ -l h_rt=01:00:00
8 species=$1
9 gffread_folder=proteins_chunks
10 gffread_chunk="$species"_proteins.fa_chunk_
11
```



```
12 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
    ec2020/step9/kofamKoala
13
14 mkdir $gffread_folder
15 module load exonerate/2.4.0
16 fastasplit -f single_fasta_proteins.fa -o ./ $gffread_folder/ -c
    100
17 cd $gffread_folder
18
19 readlink -f single_fasta_proteins.fa_* > list
20 for i in $(seq 1 100);
21 do
22     original_file=$(head -"$i" list | tail -1)
23     renamed_file="$gffread_chunk"$i
24     mv -- "$original_file" "$renamed_file"
25 done
```

Run kofamkoala on the single chunks:

KofamKoala\_universal\_chunks.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=40G
7  #$ -l h_rt=36:0:0
8  #$ -t 1-100
9  #$ -l highmem
10
11 species=$1
12 target_chunk="$species"_proteins.fa_chunk_"${SGE_TASK_ID}"
```

```
13 output_kofam="$species"_kofam_result_chunk_"${SGE_TASK_ID}".txt
14
15 echo "Working on "$species
16
17 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
ec2020/step9/kofamKoala/proteins_chunks/
18
19 module load anaconda3
20 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
3/KofamKOALA_env
21
22 exec_annotation \
23 --profile=/data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/Ko
famKOALA_env/ftp.genome.jp/pub/db/kofam/profiles/ \
24 --ko-list=/data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/Ko
famKOALA_env/ftp.genome.jp/pub/db/kofam/ko_list \
25 --cpu=6 \
26 --format=mapper \
27 --report-unannotated \
28 -o $output_kofam \
29 $target_chunk
```

Some jobs will crash because of the same error so we can divide in 100 subchunks the chunks that failed the previous jobs. With these steps we will remove those problematic subchunks from the analysis. Every subchunk removed in this way will be just the 0.01% of the total so it won't impact too much the consistency of our analyses

obtain\_rename\_chunks\_universal\_fix\_v1.sh

```
1 #!/bin/bash
2 # $ -cwd
3 # $ -pe smp 1
```

```
4  #$ -l h_vmem=1G
5  #$ -j y
6  #$ -l h_rt=01:00:00
7
8  species=$1
9  chunk=$2
10 trinotate_folder_chunk=input_trinotate_proteins_chunk_"$chunk"
11 chunk_fix=input_trinotate_proteins.fa_chunk_"$chunk"
12 trinotate_chunk=input_trinotate_proteins.fa_chunk_
13
14 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
   ec2020/step8/input_trinotate_proteins_chunks/
15
16 mkdir $trinotate_folder_chunk
17 module load exonerate/2.4.0
18 fastasplit -f $chunk_fix -o ./$trinotate_folder_chunk/ -c 100
19 cd $trinotate_folder_chunk
20
21 readlink -f "$trinotate_chunk"* > list
22 for i in $(seq 1 100);
23 do
24     original_file=$(head -"$i" list | tail -1)
25     renamed_file="$trinotate_chunk"$i
26     mv -- "$original_file" "$renamed_file"
27 done
```

Run HMMER on the subchunks

HMMER\_universal\_chunks\_fix\_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=1G
7  #$ -l h_rt=36:0:0
8  #$ -t 1-100
9  #$ -l highmem
10
11 species=$1
12 chunk=$2
13 trinotate_folder_chunk=input_trinotate_proteins_chunk_"$chunk"
14 target_chunk=input_trinotate_proteins.fa_chunk_"${SGE_TASK_ID}"
15 PFAM_out=PFAM_"${SGE_TASK_ID}".out
16 pfam_log=pfam_"${SGE_TASK_ID}".log
17
18 echo "Working on "$species" on chunk: "$chunk
19
20 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
ec2020/step8/input_trinotate_proteins_chunks/$trinotate_folder_c
hunk
21
22 module load anaconda3
23 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
3/trinotate_env
24
25 hmmscan --cpu 12 --domtblout $PFAM_out /data/SBCS-MartinDuranLab
/03-Giacomo/db/trinotate/Pfam-A.hmm $target_chunk > $pfam_log
```

to be removed:

- oasisia 38 37 rm PFAM\_37.out
- oasisia 46 25 rm PFAM\_25.out
- oasisia 75 32 rm PFAM\_32.out
- osedax 75 81 rm PFAM\_81.out
- osedax 52 70 rm PFAM\_70.out

Now we need to put together the PFAM.out file from the subchunks to basically rebuild the PFAM.out file of the chunk that failed before

merge\_chunks\_universal\_fix\_v1.sh

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -pe smp 1
4  #$ -l h_vmem=1G
5  #$ -j y
6  #$ -l h_rt=01:00:00
7
8  species=$1
9  chunk=$2
10 trinotate_folder_chunk=input_trinotate_proteins_chunk_"$chunk"
11 PFAM_chunk_out=PFAM_"$chunk".out
12
13 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
   ec2020/step8/input_trinotate_proteins_chunks/$trinotate_folder_c
   hunk
14
15 for i in $(seq 1 100);
16 do
17     PFAM_out=PFAM_"$i".out
18     PFAM_out_ok=PFAM_"$i"_ok.out
19     head -n -10 $PFAM_out | tail -n +4 > $PFAM_out_ok
20 done
21 cat PFAM*_ok.out > ../$PFAM_chunk_out
```

And finally we need to put together the PFAM.out file of all the chunks to reconstruct the PFAM.out file that we would have obtained if we didn't encounter the Segmentation fault error

## merge\_chunks\_universal\_v1.sh

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -pe smp 1
4  #$ -l h_vmem=1G
5  #$ -j y
6  #$ -l h_rt=01:00:00
7
8  species=$1
9
10 cd /data/scratch/btx654/btx604-scratch/$species/New_annotation_D
    ec2020/step8/input_trinotate_proteins_chunks
11
12 for i in $(seq 1 100);
13 do
14     PFAM_out=PFAM_"$i".out
15     PFAM_out_ok=PFAM_"$i"_ok.out
16     head -n -10 $PFAM_out | tail -n +4 > $PFAM_out_ok
17 done
18 cat PFAM*_ok.out > ../PFAM.out
```

End of fix