

# Genome annotation pipeline

By Ferdi Marlétaz, Daria Gavr, Chema Martín and Giacomo Moggioli

This pipeline is intended to be a standard procedure to annotate genomes combining different sources of evidence (RNA-seq, *de novo* transcriptomes, inter-specific protein alignments, *ab initio* predictions). At QMUL, the core elements of the pipeline are in Apocrita and are loaded into the path after sourcing a script James Crowe prepared:

```
source /data/SBCS-Informatics/chema_pipeline_0519/init.sh
```

This results in:

```
1 $ Loading trinity/2.4.0
2 $ Loading requirement: intel/2017.1 bowtie2/2.3.2
3 $ checking paths to required modules
4 $ /data/SBCS-Informatics/chema_pipeline_0519/repeatModeler/RepeatModeler-open-1.0.11/BuildDatabase
5 $ /data/SBCS-Informatics/chema_pipeline_0519/repeatModeler/RepeatModeler-open-1.0.11/RepeatModeler
6 $ /data/SBCS-Informatics/chema_pipeline_0519/repeatModeler/RepeatMasker/RepeatMasker
7 $ /share/apps/centos7/star/2.5.3a/bin/STAR
8 $ /data/SBCS-Informatics/chema_pipeline_0519/stringtie-1.3.6.Linux_x86_64/stringtie
9 $ /share/apps/centos7/trinity/2.4.0/Trinity
10 $ /data/SBCS-Informatics/chema_pipeline_0519/env/bin/mikado
```

**Note:** As an example, the following pipeline is described as it was implemented in

*Owenia*. For future projects, the name of the initial genome and downstream files should be adjusted to the species of study. To make the design of scripts easier and as further orientation in the HPC resources required in each step, I include the heading of the `.sh` files.

## Step 0: Assembly QC

Before initiating the annotation pipeline, evaluate the quality of the genome assembly following the pipeline described in [+Assembly quality check and genome size estimations](#).

## Step 1: Mask repeat elements

To avoid including transposable elements (TEs) into the gene repertoire and excluding repetitive regions that could affect mapping and gene prediction algorithms, it is necessary to mask (either with Ns [hard masking] or with lower case letters [soft masking]) the genome. To do so, we use RepeatModeler (to generate a database of TEs and repetitive sequences of the species) and RepeatMasker (to mask the genome using the previously built database). In addition to being an essential first step for gene annotation, identification of repetitive elements is a project on its own, as TEs are an essential component of animal genomes.

**Note:** RepBase (a curated database used for initial TE/repetitive element prediction requires now a licence; we have an old version [2015] in the server, but it might be good to consider buying institute-wide licence)

**Potential improvements:** considering using short-read based, assembly-free repeat prediction tools. However, as we move to PacBio-only assemblies, these might prove suboptimal.

### 1.1 RepeatModeler

The first step is to prepare a database of the unmasked genome and generate a species-specific prediction of repetitive elements. [RepeatModeler](#) uses sequence similarity approaches (with either BLAST or HMMR) to identify potential repeat elements in the genome of interest that are similar to those curated in a database (RepBase). After several rounds of search, it outputs a set of consensus repetitive sequences.

Install with conda:

```
1 module load anaconda2
2 conda create -n Repeats_env
3 source activate Repeats_env
4 conda install -c bioconda repeatmodeler
```

- repeatmodeler version 2.0.1

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/
3 #$ -o /data/scratch/btx654/
4 #$ -j y
5 #$ -pe smp 20
6 #$ -l h_vmem=20G
7 #$ -l h_rt=240:0:0
8 #$ -l highmem
9
10 echo "Working on owenia"
11
12 module load anaconda2
13 source activate Repeats_env
14
15 cd owenia/september2020/
16
17 # Build a genomic database and run RepeatModeler to generate a O
18 wenia-specific. The database is built using BLAST (ncbi flag).
19 BuildDatabase -name owenia -engine ncbi Owenia_unmasked_v082020.
20 fa
21 RepeatModeler -engine ncbi -pa 8 -database owenia
```

RepeatModeler will create a folder starting with `RM_` and will generate two main files: `consensi.fa.classified` and `families-classified.stk`. The first is the set

of consensus families predicted by RepeatModeler in fasta format. The second is the seed alignments of these families in Stockholm format. For downstream analysis, one needs the `consensi.fa.classified`, which can be renamed to include the name of the species.

## 1.2 Filtering *bona fide* genes out of the RepeatModeler library

Sometimes, RepeatModeler will interpret an expanded gene family as a repetitive element. To avoid masking out real genes (i.e. to remove false positives), one needs to check for the presence of real genes in the `consensi.fa.classified` library. To do so, one can get a curated proteome of a related species (e.g. *Capitella*, from [ENSEMBL Metazoa](#)), remove potential TEs from *Capitella*'s proteome, and then use that proteome to identify and filter out potential genes from the `consensi.fa.classified` dataset. As a reference TE database, one can use the `RepeatPeps.lib` of RepeatMasker.

```
1 #make a diamond BLAST database of RepeatPeps.lib and BLAST Capit
  ella proteome against it
2 module load diamond/0.9.22
3 diamond makedb --in RepeatPeps.lib -d RepeatPeps
4 diamond blastp -d RepeatPeps -q cte_NRproteome.fasta -o cte_NRpr
  oteome.vs.RepeatPeps.1e5.blastp -f 6 qseqid bitscore evaluate stit
  le -k 25 -e 1e-5 -p 8
5
6 #Get the unique set of qseqids with significant hit against a TE
  and filter it out of the proteome. One can do that easily with s
  hell commands (cut, grep, etc) tabulating the fasta file of the
  proteome (so that one can find at the same time the ID and seque
  nce associated to it). Alternatively, if there are not many, one
  can do that by hand.
7
8 #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.
```

```

9 diamond makedb --in Capitella_filteredProt.fasta -d CapitellaProtNoTEs
10 diamond blastp -d CapitellaProtNoTEs -q consensi.fa.classified -o consensi.fa.classified.vs.CapitellaProtNoTEs.1e10.blastp -f 6 qseqid bitscore evalue stitle -k 25 -e 1e-10 -p 8
11
12 #Since there will not be many hits, I recommend removing the potential bona fide genes out of the consensi.fa.classified by hand, as one makes sure that each of those hits are actual genes and not a TEs.

```

With *Owenia*, this approach uncovered 7 potential genes incorporated into the consensi.fa.classified database.

```

1 rnd-1_family-1155#Unknown ===== alcohol dehydrogenase
2 rnd-1_family-59#Unknown ===== histone cluster?
3 rnd-6_family-1486#Unknown ===== KR superfamily domain (plasminogen)
4 rnd-6_family-235#Unknown ===== FReD superfamily domain
5 rnd-6_family-3681#Unknown ===== amidase
6 rnd-6_family-6593#Unknown ===== adenosine receptor (!) GPCR!
7 rnd-6_family-889#DNA/Academ =====

```

I verified manually with the NCBI BLAST web server that these are real genes with clear homology to non-TEs and extracted them from the file. I saved the filtered consensi.fa.classified as *Owenia\_RMconsensusDB\_noGenes.fasta*

If you already have a fasta file of the proteome of a related species you can use this script:

First create two Condas environments with the software we will need for this step:

```

1 module load anaconda2
2 conda create -n fastx_toolkit
3 source activate fastx_toolkit
4 conda install -c bioconda fastx_toolkit

```

```
1 module load anaconda2
2 conda create -n diamond
3 source activate diamond
4 conda install -c bioconda diamond
```

Then we can run:

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/
3 #$ -o /data/scratch/btx654/
4 #$ -j y
5 #$ -pe smp 8
6 #$ -l h_vmem=15G
7 #$ -l h_rt=120:0:0
8 #$ -l highmem
9
10 capitella=/data/SBCS-MartinDuranLab/03-Giacomo/data/01-owenia/Ca
    pitella_filteredProt.fasta
11 consensi=consensi.fa.classified
12
13 cd owenia/september2020/
14 mkdir step_1.2
15 cd RM*
16 cp $consensi ../step_1.2/
17 cd ../step_1.2
18 cp $capitella ./
19
20 #convert the fasta file into a multi fasta with 60 charachters p
    er line
21 module load anaconda2
22 source activate fastx_toolkit
```

```
23 fasta_formatter -i $capitella -o Capitella_filteredProt.multi.fasta -w 60
24 source deactivate
25
26
27 #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.
28 source activate diamond
29
30 diamond makedb --in Capitella_filteredProt.multi.fasta -d CapitellaProtNoTEs
31
32 diamond blastp -d CapitellaProtNoTEs -q $consensi -o consensi.fasta.vs.CapitellaProtNoTEs.1e10.blastp -f 6 qseqid bitscore evalue stitle -k 25 -e 1e-10 -p 8 -b12 -c1
```

## 1.3 Annotate the TEs with TEclass

The annotation is made possible thanks to the online tool [TEclass](#). This software classifies unknown TEs consensus sequences into four categories, according to their mechanism of transposition: DNA transposons, LTRs, LINEs, SINEs

Feed the consensus filtered file (the file we obtained in step 1.2:

Owenia\_RMconsensusDB\_noGenes.fasta) to TEclass by uploading the whole file to the “[new request](#)” page, mind that this software is working only with fasta file weighting less than 5Mb.

This will output a txt file in fasta format containing the information about the TE type at the end of the header line, that will be used later to obtain plot about the diversity of TEs. Rename this file teclass\_owenia.txt and upload it to Apocrita.

## 1.4 RepeatMasker-Kimura

Now that we have a curated database of repetitive elements, we can run

RepeatMasker to mask those regions in the genome. To avoid mistaken Ns (gaps in the genome) from masked regions (TEs, etc), we run RepeatMasker with the `-xsmall` flag to soft-mask the genome. To generate a track of repeat elements that we can explore thereafter locally in IGV, we add the `-gff` flag.

**Note:** I run this in `Ofus_unmasked_v072019.fa`, which is the version of the assembly with the scaffolds ordered from largest to shortest.

Created a new anaconda3 env named Repeats\_env3:

```
1 module load anaconda3
2 conda create -n Repeats_env3
3 source activate Repeats_env3
4 conda install -c bioconda repeatmodeler
```

- repeatmodeler version 2.0.1

!!!!!!! IMPORTANT !!!!!!! otherwise it's not working! set some path with the script  
/data/home/btx604/.conda/envs/Repeats\_env/share/RepeatMasker/configure. I  
Chose RMBlast as default (n2) and I specified the path to the bin of the conda env:  
/data/home/btx604/.conda/envs/Repeats\_env/bin  
/data/home/btx654/.conda/envs/Repeats\_env3/share/RepeatMasker/  
/data/home/btx654/.conda/envs/Repeats\_env3/bin

The first thing we will need in this step is the genome size. To get this information one can run quast:

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

```
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=10:0:0
8
9  module load anaconda3
10 source activate quast
11
12 cd owenia/september2020/
13 mkdir quast
14 cd quast
15
16 quast \
17 /data/scratch/btx654/owenia/september2020/Owenia_unmasked_v0820
18 20.fa \
19 -o /data/scratch/btx654/owenia/september2020/quast/ \
--eukaryote \
```

With the info about the genome size we can run the following script. It should be submitted like this:

```
qsub kimura_sept2020_owenia_v1.sh owenia 500139420
```

- \$1 is the name of the species
- \$2 is the genome size in bp

`kimura_sept2020_owenia_v1.sh`

```
1  #!/bin/bash
2  # $ -wd /data/scratch/btx654/
3  # $ -o /data/scratch/btx654/
4  # $ -j y
5  # $ -pe smp 15
6  # $ -l h_vmem=15G
7  # $ -l h_rt=72:0:0
8  # $ -l highmem
9
10 genome_size=$2 #in bp
```

```
11 species_fasta=Owenia_unmasked_v082020.fa
12 species_fasta_path=/data/scratch/btx654/owenia/september2020/$species_fasta
13 species_unmasked="$1"_unmasked.fa
14 species_unmasked_align="$species_unmasked".align
15 teclass_name=teclass_"$1".txt
16 teclass_path=/data/SBCS-MartinDuranLab/03-Giacomo/data/01-owenia/sept2020/$teclass_name
17 output=TEclass_"$1"
18 divsum="$output".divsum
19 html="$output".html
20
21 echo "Working on "$1
22
23 module load anaconda3
24 source activate Repeats_env3
25
26 cd owenia/september2020/
27 mkdir kimura_highmem
28 cd RM*
29 cp consensi.fa.classified ../kimura_highmem/ #I have found no genes in step 1.2 so I am using the consensi file obtained in step 1.1
30 cd ../kimura_highmem
31 cp $species_fasta_path ./
32 mv $species_fasta $species_unmasked
33 cp $teclass_path ./
34
35 #the next section will update the filter consensi file obtained in step 1.2 with the annotations from TEclass
36 grep '>' $teclass_name | sed 's/|/:/g' | sed 's/#/:/' | awk ' B
```

```
EGIN { FS = ":" } ; { OFS="\t" ; print $1,$4}' > table_of_change
s.tsv
37 grep '>' consensi.fa.classified | sed 's/|/:/g' | sed 's/#/:/' |
awk ' BEGIN { FS = ":" } ; { OFS="\t" ; print $1}' > fasta_entr
ies
38 grep '>' $teclass_name | sed 's/|/:/g' | sed 's/#/:/' | awk ' B
EGIN { FS = ":" } ; { OFS="\t" ; print $1}' > text_entries
39
40 comm -13 <(sort text_entries) <(sort fasta_entries) > missing_en
tries
41
42 cat missing_entries | awk '{ OFS="\t" ; print $1,"unclear"}' >>t
able_of_changes.tsv
43
44 while read -r line
45 do
46 Target=$(echo $line | awk '{ print $1}')
47 rep=$(echo $line | awk '{ print $2}')
48 sed -iE 's/'"$Target"'#.* (/'"$Target"'#'"$rep"' (/ consensi.f
a.classified
49 done <table_of_changes.tsv
50 #end of the section
51
52 RepeatMasker -pa 15 -xsmall -gff -a -lib consensi.fa.classified
$species_unmasked
53
54 #Kimura distances analyses
55 /data/home/btx654/.conda/envs/Repeats_env3/share/RepeatMasker/ut
il/calcDivergenceFromAlign.pl -s $divsum $species_unmasked_align
56 /data/home/btx654/.conda/envs/Repeats_env3/share/RepeatMasker/ut
il/createRepeatLandscape_TEclass.pl -div $divsum -g $genome_siz
e > $html
```

Repeat masker generates four main outputs:

- `Ofus_unmasked_v072019.fa.masked`: the masked genome
- `Ofus_unmasked_v072019.fa.out`: a tabular collection of all masked elements, with IDs, etc
- `Ofus_unmasked_v072019.fa.tbl`: a summary table of the masking process, good for publications!
- `Ofus_unmasked_v072019.fa.out.gff`: a GFF file of the masked genome, useful to add as track in IGV.

The Kimura distances analyses will output:

- `TEclass_owenia.html`: This file is containing the info that we will need to plot Kimura values. Otherwise one can move this file to his personal pc and click the link to see the plots in a browser.

## 1.5 LTR\_Finder

RepeatModeler is a general prediction tool, and thus it might underperform while predicting specific TE families. LTR transposons are often abundant, specially in plant genomes, and there are tools that deal specifically with its prediction. To improve TE annotation, one can run structural based algorithms such as [LTR\\_Finder](#). Different from RepeatModeler, LTR\_Finder does not search for sequence similarity, but for the structural sequence elements that characterise LTR transposons.

**Note:** because RepeatModeler/RepeatMasker and LTR\_Finder follow different approaches, they can be run in parallel.

Installation of `ltr_finder` with Conda:

```
1 module load anaconda2
2 conda create -n ltr_finder
3 source activate ltr_finder
4 conda install -c bioconda ltr_finder
```

Create another conda environment that will be used to install Perl dependencies:

```
1 module load anaconda2
2 conda create -n myperl perl perl-app-cpanminus
```

- Once the environment is activated, modules can be installed via `cpan` or `cpanm` with no special configuration needed.

```
1 source activate myperl
2 cpanm My::Module
3 cpanm IO::Scalar
```

```
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=72:0:0
8
9 species_unmasked=/data/scratch/btx654/owenia/september2020/Oweni
a_unmasked_v082020.fa
10 txt_output=Owenia_Sept2020_LTRfinder_FullOutput.txt
11 gff_output=Owenia_Sept2020_LTRfinder_FullOutput.gff
12
13 cd owenia/september2020/
14 mkdir ltr_finder
15 cd ltr_finder
16
17 module load anaconda2
18 source activate ltr_finder
19
20 ltr_finder $species_unmasked -w 0 -C 2>&1 > $txt_output
21
22 conda deactivate
23 source activate myperl
24
25 perl /data/SBCS-MartinDuranLab/03-Giacomo/src/dawgpaws/MODIFIED_
```

```
cnv_ltrfinder2gff.pl -i $txt_output -o $gff_output
```

LTR\_Finder generates a `.txt` file, which needs to be converted into `.GFF` prior merging with RepeatCraft. To do so, we can use a script from the [Dawgpaw](#) package([cnv\\_ltrfinder2gff.pl](#)). I am using a modified version of it where the first line listed of being “#!/usr/bin/perl -w” is “#!/usr/bin/env perl”. When you specify env the script will look into the anaconda env for the perl modules.

And this generates `Owenia_LTRfinder_FullOutput.gff` file with all the predicted LTRs, which can be used together with the `.gff` output of RepeatMasker to run RepeatCraft and generate a curated set of repetitive elements.

## 1.6 RepeatCraft

To merge different sources of evidence and curate the potential repetitive regions of the genome, reducing artefacts, we run [RepeatCraft](#). The configuration file `Owenia.repeatcraft.cfg` looked like this, with default parameters except for `LTR_finder`, for which I calculated from the `.gff` file the largest LTR found:

calculate the longest ltr:

```
awk 'BEGIN {FS="\t"}; {print $5 == $4}' Owenia_Sept2020_LTRfinder_FullOutput.gff | sort -nr | head -n1
```

```
1 #!/bin/bash
2 #$ -wd /data/scratch/btx654/
3 #$ -o /data/scratch/btx654/
4 #$ -j y
5 #$ -pe smp 1
6 #$ -l h_vmem=6G
7 #$ -l h_rt=12:0:0
8
9 LTR_gff=Owenia_Sept2020_LTRfinder_FullOutput.gff
10 LTR_gff_original=/data/scratch/btx654/owenia/september2020/ltr_finder/$LTR_gff
```

```
11
12 RM_out="$1"_unmasked.fa.out
13 RM_out_original=/data/scratch/btx654/owenia/september2020/kimura
   _highmem/$RM_out
14
15 RM_out_gff="$1"_unmasked.fa.out.gff
16 RM_out_gff_original=/data/scratch/btx654/owenia/september2020/ki
   mura_highmem/$RM_out_gff
17
18 output="$1"_masked_RepeatCraft
19
20 cd owenia/september2020/
21 mkdir repeatcraft
22 cd repeatcraft
23
24 #create a configuration file that we can modify accordingly to t
   he longest ltr size we have found
25 echo "#configuration file for RepeatCraft
26
27 # Label short TEs
28 shortTE_size: 100
29 mapfile: None
30
31 # LTR grouping (based on LTR_FINDER result).
32 ltr_finder_gff: $LTR_gff
33 max_LTR_size: 25000
34 LTR_flanking_size: 200
35
36 # TEs grouping
37 gap_size: 150" > repeatcraft.cfg
38
```

```
39 # Load python
40 module load python
41
42 # Copy required files into the working directory
43 cp $RM_out_original ./
44 cp $RM_out_gff_original ./
45 cp $LTR_gff_original ./
46
47 # Run RepeatCraft
48 /data/SBCS-MartinDuranLab/02-Chema/src/repeatcraft/repeatcraft.
  py \
49 -r $RM_out_gff \
50 -u $RM_out \
51 -c repeatcraft.cfg \
52 -o $output
```

RepeatCraft will generate several outputs:

- `Ofus_masked_RM_RC_v072019.rclabel.gff`: a GFF file with all the evidences, renamed by RepeatCraft.
- `Ofus_masked_RM_RC_v072019.rmerge.gff`: the GFF file with all the repeats, merging all evidences.
- `Ofus_masked_RM_RC_v072019.summary.txt`: a summary of the run.

## 1.7 Soft-masking

To soft-mask the genome with the final, polished annotation of TEs and repeats, we use [BEDtools](#).

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



```
7  #$ -l h_rt=0:30:0
8
9  unmasked="$1"_unmasked.fa
10 unmasked_original=/data/scratch/btx654/owenia/september2020/kimura_highmem/$unmasked
11
12 masked="$1"_masked_RepeatCraft.rmerge.gff
13 masked_original=/data/scratch/btx654/owenia/september2020/repeatcraft/$masked
14
15 output="$1"_softmasked.fa
16
17 cd owenia/september2020/
18 mkdir softmasking
19 cd softmasking
20
21 # Load the application module
22 module load bedtools/2.26.0
23
24 # Copy required files into working directory
25 cp $unmasked_original ./
26 cp $masked_original ./
27
28 bedtools maskfasta -soft -fi $unmasked -bed $masked -fo $output
```

This generates the soft-masked genome assembly `Owenia_softmasked_v072019.fa`, which will be used in all downstream analyses.

## 1.8 Kimura plots

In the html file obtained during the Kimura distances analyses in step 1.4 we can find the informations required to plot these data with R.

First we need to copy-paste these info in order to generate two files that look like this:

table:

	Divergence	unclear	nonLTR	DNA	LTR	LINE	Retro	SINE	
1	0,	0.206642179894558,	0.00297577023622733,	1.61080164406957,	0.3	16742279582761,	0.494263179654985,	0.149770238066817,	0.00494202
2	196659483								
3	1,	0.29488937304722,	0.00398628846332489,	2.82234021865343,	0.26	9521446639819,	0.574690153397627,	0.271607864862962,	0.005443482
4	13944024								
5	2,	0.317798585042547,	0.0040632669986301,	2.74432437259195,	0.26	9168145154405,	0.484371937728884,	0.226616010391662,	0.004372180
6	86108869								
7	3,	0.397777483726438,	0.00280301840634757,	2.45014800073148,	0.2	04195662081585,	0.477379087615209,	0.186876091470654,	0.00636262
8	584540927								
9	4,	0.355636834225145,	0.00214020322573254,	2.20615643533957,	0.1	48996653773062,	0.3568137060662,	0.168448429839823,	0.0103765066
10	148955								
11	5,	0.301087644721146,	0.00204822887186137,	1.97399417146523,	0.1	44916791401885,	0.302978117581694,	0.149525706252069,	0.00934619
12	390729089								
13	6,	0.262411229252835,	0.00260507360127702,	1.79007305602906,	0.1	38096093285348,	0.299755616143994,	0.124288943271058,	0.00953194
14	211326114								
15	7,	0.239363255949711,	0.00186527988535677,	1.48807726453556,	0.1	22190328448815,	0.344873835379743,	0.112410055580102,	0.00886492
16	810344764								
17	8,	0.242699125775769,	0.00108749676240277,	1.37310192425944,	0.1	23240835525422,	0.313344626984212,	0.116041243059785,	0.00878874
18	93451326								
19	9,	0.193704987301341,	0.000526653148036202,	1.25317876363355,	0.	119865776626845,	0.232635132019788,	0.108142645504727,	0.0078788
20	0307455069								
21	10,	0.177116012970943,	0.000409885707469329,	1.22066243048788,					

0.0983861660014722, 0.178551012835581, 0.114094785809925, 0.0104  
884753935213

13 11, 0.154375753864792, 0.000486464354279453, 1.09288665948387,  
0.0964435076923151, 0.164249200752862, 0.104452474472018, 0.0123  
475570072041

14 12, 0.137114766918392, 0.000577039098417797, 1.04859281038075,  
0.0893678806601567, 0.158548990199573, 0.097118719416278, 0.0122  
175932462992

15 13, 0.134265961279357, 0.000494462124181293, 0.924270476420355,  
0.0887288588449997, 0.164660686014312, 0.0890695638428181, 0.012  
7408473421271

16 14, 0.162618455469877, 0.000565642276307674, 0.854410756104768,  
0.0942173284401377, 0.164392960666848, 0.0812147540779729, 0.013  
3312827051305

17 15, 0.199831878878893, 0.000802176321154609, 0.778722061140472,  
0.0925769858332702, 0.152439493771557, 0.0804991536160057, 0.011  
7717175742716

18 16, 0.20533514434835, 0.000530452088739576, 0.683648771376589,  
0.0843664752520407, 0.132561836457522, 0.0759090335250919, 0.010  
7180113897041

19 17, 0.160126350368463, 0.000905147608640807, 0.616042022842351,  
0.0803999812692229, 0.123426183842897, 0.0701362432099433, 0.008  
51922449944058

20 18, 0.118752687000757, 0.000533451252452766, 0.577206851641488,  
0.0727209224979707, 0.122502641363482, 0.0641169216375706, 0.008  
58040743918966

21 19, 0.108250615398402, 0.000453273609186814, 0.538037773547224,  
0.0721268881385115, 0.121480726314275, 0.0585328786921055, 0.007  
14660724003719

22 20, 0.0852244360182607, 0.000310313472191414, 0.51577038258652,  
0.0651308389168764, 0.124223361557863, 0.0546997475224009, 0.006  
76931244491786

23 21, 0.0644452300920411, 0.00027992194656442, 0.464038007641949,  
0.0648825081614243, 0.124558268172503, 0.0533285298727303, 0.006  
47319501430221

24 22, 0.0533409264160781, 0.000654417522218105, 0.405780052290219,  
0.0639939559253298, 0.121314572644564, 0.0490565210796621, 0.005  
63922755778779

25 23, 0.0496717495293612, 0.000647219529306448, 0.366113912796556,  
0.0583603268064733, 0.105190068801215, 0.0476709074441683, 0.006  
11229564748166

26 24, 0.048097788412679, 0.000496861455151845, 0.323490797825934,  
0.0501552147199275, 0.0905249580207055, 0.0447853120635842, 0.00  
424301687717397

27 25, 0.040797823934774, 0.000382093457060433, 0.284645629412695,  
0.0431381713523001, 0.0969303719350896, 0.0455602959670725, 0.00  
32454950261669

28 26, 0.0388883563707096, 0.000302315702289574, 0.260197046655511,  
0.0384094898978369, 0.129244161557991, 0.0347936981252148, 0.002  
88619521332672

29 27, 0.0302935529456966, 0.000119966548527609, 0.228435702988579,  
0.035218179762755, 0.128233843275141, 0.0304874988658163, 0.0027  
9761991166383

30 28, 0.0267311462871693, 7.35794830969332e-05, 0.194351407053657,  
0.0289337321181362, 0.0829110810741533, 0.0294119987582662, 0.00  
14673908327402

31 29, 0.0257762125608895, 5.69841105506141e-05, 0.156749691915906,  
0.0276886792886671, 0.0533801154885972, 0.0312020996065457, 0.00  
108149843497639

32 30, 0.0220198599822426, 8.99749113957064e-05, 0.129004028516688,  
0.0224695345949735, 0.0615684322583491, 0.0266597661907954, 0.00  
0780382398172094

33 31, 0.0188185526347833, 1.05970451199388e-05, 0.11316144606238,  
0.0194077883323014, 0.0617273879351482, 0.0184102664812944, 0.00  
116727451717363

34 32, 0.0163428429616686, 2.57928079334358e-05, 0.095698715370206  
2, 0.0211980891248284, 0.0641169216375706, 0.0257248268892702,  
0.000736994496454609

35 33, 0.0161700911317888, 0, 0.0754889506609977, 0.022795643662721  
1, 0.0611481494500074, 0.0219538783805524, 0.000502060005588042

36 34, 0.0141248614236406, 7.39793715920253e-06, 0.066701201037102  
8, 0.0284214749559233, 0.0613606901851488, 0.0177364543670643,  
0.00047306808969387

37 35, 0.0118005095459182, 0, 0.0488969655701204, 0.020502683031863  
4, 0.0691667135535927, 0.0197624894274481, 0.000333307060659206

38 36, 0.00975527983777004, 0, 0.0396305494176004, 0.02368819478376  
65, 0.0555577082886208, 0.0172511896782701, 0.000330707785441108

39 37, 0.0103531131379326, 0, 0.0314224381673414, 0.022952200008549  
6, 0.0522682255279938, 0.0136315989649446, 8.83753574153383e-05

40 38, 0.00961771819545838, 0, 0.0270224650558438, 0.02154779161378  
64, 0.0512495095867468, 0.0111178998847961, 0.000179749878543867

41 39, 0.00990303863670654, 0, 0.0187461728171717, 0.02296519638464  
01, 0.0396107549370933, 0.0109699411416121, 0.000122565823745707

42 40, 0.00724497980982983, 0, 0.0157386114455845, 0.02075301322979  
1, 0.0354063272996958, 0.00981366355805347, 2.69924734187119e-05

43 41, 0.00501500161694913, 0, 0.0139589076981774, 0.02927063817525  
12, 0.0297421067109647, 0.00934199507809243, 0

44 42, 0.00417603555424605, 0, 0.0131005470434624, 0.03654600951070  
8, 0.0271676245795622, 0.0061654808173289, 2.21938114776076e-05

45 43, 0.00439857350176477, 0, 0.00904947664393261, 0.0228564267139  
751, 0.020711624770549, 0.00572680313821294, 0

46 44, 0.00258487923227487, 0, 0.00997721795254611, 0.0198278711963  
956, 0.0133648733387182, 0.0036017956752939, 0

47 45, 0.00342124601975985, 0, 0.00613129035099853, 0.0148008729245  
937, 0.00786740625244057, 0.00300676159459696, 0

48 46, 0.00157935961136597, 0, 0.00429420260454575, 0.0097076931068  
5408, 0.00422702133737029, 0.00252429612526843, 0

49 47, 0.000431879574699391, 0, 0.00419842930997121, 0.004476151869  
81262, 0.00266225765607518, 0.00116087630125216, 0

50 48, 0.00114268137472547, 0, 0.00105030713235921, 0.0042124254072  
9943, 0.00183668785795769, 0.00064561997532608, 0

51 49, 0.00030451508901258, 0, 0.000792578997272401, 0.002318753438  
79113, 0.00115787713753897, 0.00086335926090369, 0

52 50, 0.000339305388085586, 0, 0.000576239321427613, 0.00093973796

```
3466267, 0.000646019863821172, 7.51790370773014e-05, 0
```

## piechart

```
1 Simple_repeat 3121129
2 SINE 856876
3 Retro 14621610
4 LINE 32471055
5 LTR 16839526
6 DNA 147110666
7 nonLTR 148057
8 unclear 21891503
9 Unmasked 263078998
```

Then we have to modify these two files with the commands:

```
sed -e 's/ /\t/g' table | sed -e 's/[,,//g' > table_owenia_sept2020
```

to obtain the data for the piechart with percentage

```
1 genome_size=500139420 #obtained with quast in step 1.4
2 awk '{ print $1, 100*$2/500139420 }' pie > pie_owenia_sept2020
```

- then we should manually add a line containing “name percentage” as the first line of the file pie\_owenia\_sept2020

We can finally plot our data using Rstudio. The datasets should be imported from text files

```
1 library(ggplot2)
2 library(plyr)
3 library(reshape2)
4 library("RColorBrewer")
5
6 #this will generate the Kimura plot using pie_owenia_sept2020
7 lm <- melt(table_owenia_sept2020, id.vars=0:1)
```

```

8
9  ggplot(data=lm, aes(x=lm$Div, y=lm$value, fill=lm$variable))+geom
  m_bar(stat="identity")+labs(x="Kimura Substitution Level (%)", y
    ="Genome Proportion")+ scale_fill_brewer(name=element_blank(), p
      alette="Set3")+theme_bw(base_size = 14)
10
11  #this will generate the piechart using pie_owenia_sept2020
12  ggplot(pie_owenia_sept2020, aes(x = "", y=percentage, fill=name,
    order=pie_owenia_sept2020$name)) +
13      geom_bar(width = 1, stat = "identity", color = "white") +
14      coord_polar("y", start = 0) +
15      scale_fill_brewer (name=element_blank(), palette="Set3") +
16      theme_void()

```

- The plots have been exported in pdf 5x6.5 inch

## Step 2: Generate gene evidences

This step will align different sources of evidence (e.g. short Illumina reads from RNA-seq data, a *de novo* transcriptome assembly generated with Trinity, spliced-alignments of proteomes of closely-related species, etc) to the soft-masked genome assembly. Each source of evidence can be process independently and in parallel, so that jobs run simultaneously and this step is done faster.

### 2.1 *De novo* transcriptome assembly with Trinity

Independently assembled transcript evidences will be the base for genome annotation. For *Owenia*, we use three different *de novo* transcriptomes:

- Harald Hausen's *Owenia* reference transcriptome (the one in the BLAST server). Already assembled!
- Stage-specific RNA-seq data from oocyte to larva (Illumina sequenced)
- Tissue-specific RNA-seq data from the adult (BGI sequenced)

#### 2.1.1 *De novo* assembly of Illumina strand-specific embryonic data

```

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
```

Rename the `trinity.fasta` file to `Oxford.trinity.fasta` and run transcriptome stats script:

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -pe smp 1
5  #$ -l h_vmem=1G
6  #$ -l h_rt=1:0:0
7
8  module load perl
9  /share/apps/trinity/r20140717/util/TrinityStats.pl Oxford_Illumi
    na_trinity_R1/Trinity.fasta
```

```
1  Total trinity 'genes': 516585
2  Total trinity transcripts: 1111027
3  Percent GC: 37.16
4  Contig N50: 620
```

### 2.1.2 *De novo* assembly of BGI strand-specific adult data

BGI provides cleaned data (adaptors removed and quality trimmed reads).

Sequencing is based on a rolling circle approach, and paired reads are in a RF orientation. Because we start from several libraries, each for a different tissue, it is better to first define a sample text file (`tissueLibraries.txt`). In the file, the first column is the tissue, then the replicate (in this case there is only one replicate), then the read 1 and the read 2. Columns are tabulated (separated by tabs)

```
1  BV      BV_R1    Blood-vessel_1.fq.gz    Blood-vessel_2.fq.gz
```

2	BW	BW_R1	Body-Wall_1.fq.gz	Body-Wall_2.fq.gz
3	G	G_R1	Gut_1.fq.gz	Gut_2.fq.gz
4	H	H_R1	Head_1.fq.gz	Head_2.fq.gz
5	HC	HC_R1	Head+Chaetae_1.fq.gz	Head+Chaetae_2.fq.gz
6	O	O_R1	Ovary_1.fq.gz	Ovary_2.fq.gz
7	R	R_R1	Retractor-muscle_1.fq.gz	Retractor-muscle_2.fq.gz
8	T	T_R1	Tail_1.fq.gz	Tail_2.fq.gz
9	Ts	Ts_R1	Testis_1.fq.gz	Testis_2.fq.gz

We can launch Trinity with the following command. **Important!!** Because of the way BGI sequences (DNA nanoballs and rolling circle replication), the orientation of pair-end reads is FR, differently from the strand-specific system of Illumina, which is in RF orientation.

```

1  #!/bin/bash
2  #$ -pe smp 20
3  #$ -l h_vmem=3G
4  #$ -l h_rt=120:0:0
5  #$ -cwd
6  #$ -j y
7
8  module load trinity/2.4.0
9
10 Trinity --seqType fq --max_memory 60G --samples_file tissueLibra
    ries.txt --SS_lib_type FR --CPU 20 --output BGI_tissues_trinity

```

Rename the `trinity.fasta` file to `BGI.trinity.fasta` and run transcriptome stats script:

```

1  #!/bin/bash

```

```
2  #$ -cwd
3  #$ -j y
4  #$ -pe smp 1
5  #$ -l h_vmem=1G
6  #$ -l h_rt=1:0:0
7
8  module load perl
9  /share/apps/trinity/r20140717/util/TrinityStats.pl BGI.trinity.f
   asta
```

```
1  Total trinity 'genes':  213541
2  Total trinity transcripts:      574180
3  Percent GC: 36.74
4  Contig N50: 687
```

## 2.2 Alignment of *de novo* transcriptomes to genome with GMAP

In those cases in which there is a pre-existing *de novo* transcriptome available (e.g. assembled with Trinity, as it is the case with *Owenia*), one can directly generate spliced-alignments of the transcripts to the genome with [GMAP](#). As an example for the BGI assembly:

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -pe smp 5
5  #$ -l h_vmem=10G
6  #$ -l h_rt=10:0:0
7
8  # Activate the genome annotation pipeline environment installed
   by James and add all modules into the $PATH
9  source /data/SBCS-Informatics/chema_pipeline_0519/init.sh
```

```

10
11 # Bring input files to working directory
12 cp /data/SBCS-MartinDuranLab/02-Chema/00-OweniaGenome/00-DATA/Ow
   enia_softmasked_v072019.fa ./
13
14 # Make a GMAP index and align transcripts to the genome
15 gmap_build -D /data/scratch/btx333/07-GMAP/00-Trinity_BGI/gmapdb
   -d Owenia_softmasked Owenia_softmasked_v072019.fa
16
17 gmap -D /data/scratch/btx333/07-GMAP/00-Trinity_BGI/gmapdb -d Ow
   enia_softmasked -f 3 -n 0 -x 50 -t 5 -B 4 --gff3-add-separators=
   0 ./BGI_tissues_trinity/Trinity.fasta > BGIissuesVsGenome_gmap.
   gff3

```

The output is `BGIissuesVsGenome_gmap.gff3`, a GFF file with the aligned transcripts (94.03% alignment rate). The statistics of the aligning process are saved as `stdout` in the associated job file. Repeat the same with the Oxford-based assembly and the existing transcriptome generated by Harald Hausen (output `HaraldRefTransVsGenome_gmap.gff3`).

Transcriptome	No queries	Unaligned (no paths found)	Alignment rate
BGI adult tissues	574,180	34,497	93.99%
Oxford embryos (R1)	1,111,027	111,172	89.99%
Ofus_RefTrans (Sars)	226,138	8,837	96.09%

## 2.3 Illumina short read alignments with STAR

Simultaneously, we can align Illumina short reads (paired end, PE and single end, SE) to the soft masked genome employing [STAR](#). To avoid generating a very large single BAM file that will require an enormous amount of RAM (thus either stopping

the job or making a very long queuing), it is better to align each library alone. We can do that with an array job. Thereafter, we can use [StringTie](#) to generate a GTF file with the assembled transcripts created by STAR.

### 2.3.1 Downloading additional libraries from SRA

In order to maximise the amount of data for genome annotation, it is advisable to check at the SRA for already published Illumina libraries of the species of interest. In the case of *Owenia*, there is the dataset that was published in the annelid phylogeny. To download, generate a `.fastq` file of the reads, and trim the reads, one can run a code like so:

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -pe smp 4
5  #$ -l h_vmem=4G
6  #$ -l h_rt=12:0:0
7
8  # Load sratools
9  module load sratools/2.8.2-1
10 module load trimmomatic/0.36
11
12 # Download SRA dataset, convert it to fastq and run trimmomatic
13 prefetch SRR1222288
14
15 cp /data/home/btx333/ncbi/public/sra/SRR1222288.sra /data/scratch/btx333/08-STAR
16 cd /data/scratch/btx333/08-STAR
17 fastq-dump --gzip --split-files SRR1222288.sra
18
19 java -jar /share/apps/trimmomatic/0.36/trimmomatic-0.36.jar PE -
  threads 1 -phred33 -trimlog trimSRR1222288 ./SRR1222288_1.fastq.
  gz ./SRR1222288_2.fastq.gz SRR1222288_1_paired.fastq.gz SRR12222
```

```
88_1_unpaired.fastq.gz SRR1222288_2_paired.fastq.gz SRR1222288_2
_unpaired.fastq.gz ILLUMINACLIP:/share/apps/trimmomatic/0.36/adap
ters/TruSeq3-PE-2.fa:2:30:10 LEADING:28 TRAILING:28 SLIDINGWIND
OW:4:15 MAXINFO:40:0.5 MINLEN:36 > SRR1222288.log 2>&1
```

**Note2:** before aligning the libraries identify whether they are strand-specific or not. This is an important value that has to be specified at certain steps. For Oceane's RNA-seq, they are not stranded.

### 2.3.2 Mapping short read libraries to the genome

First generate a genome STAR index. Since this will be common to all libraries, I do not include it in the array job.

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -pe smp 4
5  #$ -l h_vmem=4G
6  #$ -l h_rt=1:0:0
7
8  # Activate the genome annotation pipeline environment
9  source /data/SBCS-Informatics/chema_pipeline_0519/init.sh
10
11 # Bring input files into working directories
12 cp /data/SBCS-MartinDuranLab/02-Chema/00-OweniaGenome/00-DATA/Ow
   enia_softmasked_v072019.fa ./
13
14 # Make genome directory
15 mkdir ./Owenia_genome_STAR
16
17 # Build STAR genome index and align the reads generating a BAM s
   orted by coordinate and a stranded wiggle file to visualise the
   RNA-seq
```

```
18 STAR --runMode genomeGenerate --runThreadN 5 --genomeDir ./Owenia_genome_STAR --genomeFastaFiles ./Owenia_softmasked_v072019.fa
```

Then generate an array job that maps every PE library to *Owenia*'s genome and transforms the BAM into a GTF. To simplify the array job, I rename all libraries to consecutive numbers:

1	Blood-vessel	1
2	Body-Wall	2
3	Gut	3
4	Head	4
5	Head+Chaetae	5
6	Ovary	6
7	Retractor-muscle	7
8	Tail	8
9	Testis	9
10	owe_ace_R1	10
11	owe_ace_R2	11
12	owe_1cell_R1	12
13	owe_1cell_R2	13
14	owe_2cell_R1	14
15	owe_2cell_R2	15
16	owe_4cell_R1	16
17	owe_4cell_R2	17
18	owe_8cell_R1	18
19	owe_8cell_R2	19
20	owe_3h_R1	20
21	owe_3h_R2	21
22	owe_4h_R1	22
23	owe_4h_R2	23
24	owe_5h_R1	24

25	owe_5h_R2	25
26	owe_9h_R1	26
27	owe_9h_R2	27
28	owe_13h_R1	28
29	owe_13h_R2	29
30	owe_18h_R1	30
31	owe_18h_R2	31
32	owe_27h_R1	32
33	owe_27h_R2	33
34	SRR1222288	34

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -pe smp 5
4  #$ -l h_vmem=5G
5  #$ -j y
6  #$ -l h_rt=8:00:00
7  #$ -t 1-9
8
9  # Activate the genome annotation pipeline environment
10 source /data/SBCS-Informatics/chema_pipeline_0519/init.sh
11
12 # Build STAR genome index and align the reads generating a BAM s
    orted by coordinate, and a stranded wiggle file to visualise the
    RNA-seq
13 STAR --runThreadN 5 --genomeDir /data/scratch/btx333/00-BGI_mapp
    ing/Owenia_genome_STAR --readFilesCommand zcat --readFilesIn ${S
    GE_TASK_ID}_1.fq.gz ${SGE_TASK_ID}_2.fq.gz --outSAMtype BAM Sort
    edByCoordinate --outWigType wiggle --outWigStrand Stranded --out
    FileNamePrefix ./BGI_library_${SGE_TASK_ID}_
14
```



```

15 # Assemble a GTF file from the BAM file
16 stringtie ./BGI_library_${SGE_TASK_ID}_Aligned.sortedByCoord.out.bam -p 5 -o BGI_library_${SGE_TASK_ID}_transcripts.gtf -v

```

The output is a series of `BGI/Oxford_library_<n>_transcripts.gtf` files, together with the sorted `BAM` files and `wiggle` for each strand, which can be uploaded to the genome browser to visualise the RNA-seq data. Before uploading the `wiggle` files, re-convert the naming to the specific library (e.g. library 9 back to testis, etc)

Library	No reads	Uniquely mapped reads %	% of reads unmapped
Blood-vessel	26,229,562	69.63	22.65
Body-Wall	25,030,358	83.16	12.01
Gut	26,223,150	84.17	11.94
Head	25,135,465	80.84	14.5
Head+Chaetae	25,164,246	86.52	9.15
Ovary	25,127,050	79.03	17.01
Retractor-muscle	26,221,434	85.39	6.95
Tail	25,090,293	79.08	16.14
Testis	26,332,262	91.92	2.62
owe_ace_R1	42,760,709	70.64	24.92
owe_ace_R2	41,755,797	67.93	27.6
owe_1cell_R1	41,650,434	72.11	22.44
owe_1cell_R2	47,164,566	74.43	20.8
owe_2cell_R1	44,590,774	71.93	22.85
owe_2cell_R2	43,188,381	72.87	22.36
owe_4cell_R1	45,902,725	72.74	22.38
owe_4cell_R2	49,945,828	72.58	22.49
owe_8cell_R1	45,680,904	72.80	22.08

owe_8cell_R2	49,401,610	73.74	21.57
owe_3h_R1	47,976,852	69.97	24.9
owe_3h_R2	44,035,927	70.40	24.65
owe_4h_R1	45,997,615	71.61	23.27
owe_4h_R2	52,692,310	73.76	21.41
owe_5h_R1	54,631,086	71.55	23.07
owe_5h_R2	46,947,697	74.61	20.45
owe_9h_R1	52,833,442	67.74	26.35
owe_9h_R2	49,640,037	74.22	20.34
owe_13h_R1	51,921,579	72.55	21.91
owe_13h_R2	51,959,906	72.48	21.95
owe_18h_R1	41,396,469	70.46	24.15
owe_18h_R2	49,610,840	76.32	18.26
owe_27h_R1	41,087,129	75.59	18.27
owe_27h_R2	59,243,774	76.69	18.13
SRR1222288	23,309,147	75.57	18.44

### 2.3.1 Convert Wig tracks to BigWig

`wig` tracks are a useful way to visualise RNA-seq data, but occupy quite a lot of space. Therefore, it is better to convert them to BigWig format, which is a binary format. We can use the [UCSC tool](#) to do so. For BGI libraries:

```

1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -l highmem
5  #$ -pe smp 4
6  #$ -l h_vmem=5G
7  #$ -l h_rt=12:0:0
8  #$ -t 1-34

```

```
9
10 wigToBigWig BGI_library_${SGE_TASK_ID}_Signal.UniqueMultiple.str
    1.out.wig ChromSizes BGI_library_${SGE_TASK_ID}_Signal.UniqueMul
    tiple.str1.out.bw
11
12 wigToBigWig BGI_library_${SGE_TASK_ID}_Signal.UniqueMultiple.str
    2.out.wig ChromSizes BGI_library_${SGE_TASK_ID}_Signal.UniqueMul
    tiple.str2.out.wig
```

## 2.4 Generate curated intron junctions with Portcullis

In order to improve prediction of intron junctions, it is important to generate a well-curated dataset from the RNA-seq alignments. We can do so with [Portcullis](#), which can be easily installed via `conda`.

```
1 module load anaconda2
2 conda config --add channels bioconda
3 conda create --yes --name portcullis portcullis
```

Once created, one can load `conda`, and activate the environment with `portcullis` to run the pipeline. It is better (and faster) to run `portcullis` on each `BAM` file independently and thereafter merge the junction `BED` files.

```
1 #!/bin/bash
2 #$ -cwd
3 #$ -j y
4 #$ -l highmem
5 #$ -pe smp 4
6 #$ -l h_vmem=5G
7 #$ -l h_rt=12:0:0
8 #$ -t 1-34
9
10 module load anaconda2
```

```
11 source activate portcullis
12
13 #Run portcullis
14 portcullis full -t 4 -v --bam_filter --orientation FR --stranded
    ness firststrand -o portcullis_library_${SGE_TASK_ID} Owenia_sof
    tmasked_v082020.fa Library_${SGE_TASK_ID}_Aligned.sortedByCoord.
    out.bam
15
16 #Deactivate environment
17 source deactivate
```

Portcullis generates multiple files (e.g. a filtered BAM), but also the curated set of junctions in `portcullis_out/3-filt/portcullis_filtered.pass.junctions.bed`, which are the ones we will use in downstream analyses ([Step 3](#)).

To merge the `bed` files, we can use `junctools`, which comes with `portcullis`:

```
1 #!/bin/bash
2 #$ -cwd
3 #$ -j y
4 #$ -l highmem
5 #$ -pe smp 1
6 #$ -l h_vmem=20G
7 #$ -l h_rt=12:0:0
8
9 module load anaconda2
10 source activate portcullis
11
12 junctools set -m 1 --operator max -o Owenia_junctions_consensus.
    bed consensus portcullis_library_1/3-filt/portcullis_filtered.pa
    ss.junctions.bed portcullis_library_2/3-filt/portcullis_filtere
    d.pass.junctions.bed portcullis_library_3/3-filt/portcullis_filt
    ered.pass.junctions.bed portcullis_library_4/3-filt/portcullis_f
```

```
filtered.pass.junctions.bed portcullis_library_5/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_6/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_7/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_8/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_9/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_10/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_11/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_12/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_13/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_14/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_15/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_16/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_17/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_18/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_19/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_20/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_21/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_22/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_23/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_24/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_25/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_26/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_27/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_28/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_29/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_30/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_31/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_32/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_33/3-filt/portcullis_filtered.pass.junctions.bed portcullis_library_34/3-filt/portcullis_filtered.pass.junctions.bed
```

13

14 source deactivate

This generates 749,824 distinct junctions.

## Step 3: Merging all gene evidences with Mikado

**Mikado** is a pipeline that allows generating the “best” set of transcripts from multiple transcript assemblies. Employing multiple input files and evidences (e.g. RNA-seq alignments, BLAST homology, curated splice junctions, etc), Mikado will select the best-scoring transcripts as the primary transcript of their respective gene loci. The final output will be a curated transcriptome-based genomic annotation that can thereafter be used as experimental hints for gene prediction ([Step 4](#)).

**Note:** It is important to use the last version of Mikado, as it has significant improvements in performance. James has installed an environment (see [below](#)).

### 3.1 Install Mikado in your user

To avoid conflicting environments with the Mikado copy installed in the initial bundle, install the newest Mikado starting from a clean shell. Then:

```
1 module load anaconda3
2 cd /data/SBCS-Informatics/chema_pipeline_0519/mikado
3 conda env create -f environment.yml
4 source activate mikado2.1
5 pip install dist/*whl
```

This should have installed the environment, which can be called with `conda activate mikado2.1`, and removed with `conda deactivate`. However, to make Mikado run, you need to copy the scoring file into a folder that Mikado can access (the home folder will change from user to user):

```
cp human.yaml /data/home/btx333/.conda/envs/mikado2/lib/python3.7/site-packages/Mikado/configuration/scoring_files/
```

After doing that, Mikado should be ready to go.

### 3.2 Generate configuration file and merge all input files into a single “all-contained” GTF file

As a first step, one needs to generate a configuration file that tells Mikado where all

input evidences are, what type they are, whether they are stranded or not, and what value should Mikado add to each of them. For Owenia, that

Owenia\_mikadoInput\_list file looks like so:

```
1 /data/scratch/btx333/10-Mikado/mikado_input/BGI_library_1_transc
ripts.gtf      stringtie1      True      0.5
2 /data/scratch/btx333/10-Mikado/mikado_input/BGI_library_2_transc
ripts.gtf      stringtie2      True      0.5
3 /data/scratch/btx333/10-Mikado/mikado_input/BGI_library_3_transc
ripts.gtf      stringtie3      True      0.5
4 /data/scratch/btx333/10-Mikado/mikado_input/BGI_library_4_transc
ripts.gtf      stringtie4      True      0.5
5 /data/scratch/btx333/10-Mikado/mikado_input/BGI_library_5_transc
ripts.gtf      stringtie5      True      0.5
6 /data/scratch/btx333/10-Mikado/mikado_input/BGI_library_6_transc
ripts.gtf      stringtie6      True      0.5
7 /data/scratch/btx333/10-Mikado/mikado_input/BGI_library_7_transc
ripts.gtf      stringtie7      True      0.5
8 /data/scratch/btx333/10-Mikado/mikado_input/BGI_library_8_transc
ripts.gtf      stringtie8      True      0.5
9 /data/scratch/btx333/10-Mikado/mikado_input/BGI_library_9_transc
ripts.gtf      stringtie9      True      0.5
10 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_10_tr
anscripts.gtf  stringtie10     True      0.5
11 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_11_tr
anscripts.gtf  stringtie11     True      0.5
12 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_12_tr
anscripts.gtf  stringtie12     True      0.5
13 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_13_tr
anscripts.gtf  stringtie13     True      0.5
14 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_14_tr
anscripts.gtf  stringtie14     True      0.5
15 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_15_tr
anscripts.gtf  stringtie15     True      0.5
```

```
16 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_16_tr
   anscripts.gtf  stringtie16      True    0.5
17 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_17_tr
   anscripts.gtf  stringtie17      True    0.5
18 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_18_tr
   anscripts.gtf  stringtie18      True    0.5
19 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_19_tr
   anscripts.gtf  stringtie19      True    0.5
20 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_20_tr
   anscripts.gtf  stringtie20      True    0.5
21 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_21_tr
   anscripts.gtf  stringtie21      True    0.5
22 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_22_tr
   anscripts.gtf  stringtie22      True    0.5
23 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_23_tr
   anscripts.gtf  stringtie23      True    0.5
24 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_24_tr
   anscripts.gtf  stringtie24      True    0.5
25 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_25_tr
   anscripts.gtf  stringtie25      True    0.5
26 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_26_tr
   anscripts.gtf  stringtie26      True    0.5
27 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_27_tr
   anscripts.gtf  stringtie27      True    0.5
28 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_28_tr
   anscripts.gtf  stringtie28      True    0.5
29 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_29_tr
   anscripts.gtf  stringtie29      True    0.5
30 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_30_tr
   anscripts.gtf  stringtie30      True    0.5
31 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_31_tr
   anscripts.gtf  stringtie31      True    0.5
```



```

32 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_32_tr
   anscripts.gtf    stringtie32    True    0.5
33 /data/scratch/btx333/10-Mikado/mikado_input/Oxford_library_33_tr
   anscripts.gtf    stringtie33    True    0.5
34 /data/scratch/btx333/10-Mikado/mikado_input/SRA_library_34_trans
   cripts.gtf       stringtie34    True    0.5
35 /data/scratch/btx333/10-Mikado/mikado_input/BGItissuesVsGenome_g
   map.gff3         gmap1     False   -0.5
36 /data/scratch/btx333/10-Mikado/mikado_input/HaraldRefTransVsGeno
   me_gmap.gff3     gmap2     False   -0.5
37 /data/scratch/btx333/10-Mikado/mikado_input/OxfordtissuesR1VsGen
   ome_gmap.gff3    gmap3     False   -0.5

```

To generate protein homology, one also needs to generate a BLAST database. Download the curated [Swissprot](#) database and upload it to Apocrita (to `/data/SBCS-MartinDuranLab/00-BlastDBs/`). Rename it with the date of download (these databases change regularly, and it is important to re-download it regularly). Use DIAMOND Blast to generate a BLAST database.

```

1 module load diamond/0.9.22
2 mv uniprot_sprot.fasta Uniprot_SwissProt_160719.fasta
3 diamond makedb --in Uniprot_SwissProt_160719.fasta --db Uniprot_
  SwissProt_160719

```

Then, run Mikado to generate a configuration file that tells the program how genes should look like (for that we initially use the score parameters based on human/vertebrate genes, but for future annotation versions and other annelid genomes, we should consider adjusting these scores to annelid genomes, so that Mikado is as accurate as possible). Once the configuration file is done, Mikado will collect all evidences and merge them in a single annotation (that is `mikado prepare` )

```

1 #!/bin/bash
2 # $ -cwd

```

```
3  # $ -j y
4  # $ -pe smp 6
5  # $ -l h_vmem=3G
6  # $ -l h_rt=24:0:0
7
8  # Activate the genome annotation pipeline environment installed
  # by James and add all modules into the $PATH
9  module load anaconda3
10 conda activate mikado2.1
11
12 # Bring genome into the working directory
13 cp /data/SBCS-MartinDuranLab/02-Chema/00-OweniaGenome/00-DATA/Ow
  enia_softmasked_v082020.fa ./
14
15 # All required files (gtfs and filtered junctions) are in the fo
  lder ./mikado_input and the file Owenia_mikadoInput_list.txt poi
  nts to them. First step is generate configuration file
16 mikado configure -t 6 --list Owenia_mikadoInput_list.txt --refer
  ence Owenia_softmasked_v082020.fa --mode permissive --scoring ma
  mmalian.yaml --copy-scoring mammalian.yaml --junctions ../00-Mik
  ado_input/Owenia_junctions_consensus.bed -bt /data/SBCS-MartinDu
  ranLab/00-BlastDBs/Uniprot_SwissProt_160719.fasta configuration.
  yaml
```

Before running the next step, open the `configuration.yaml` file and modify the parameter `max_intron_length` for the `prepare` section. By default, it is set to 1,000,000 bp (1Mb), but that generates problems at the last `pick` stage, because Mikado generates super locus that are too complex to resolve. Considering the *Capitella* and *Dimorphilus* genome, a reasonable maximum intron length is 50kb (50,000). Once that is changed, we can run `mikado prepare`:

```
1  #!/bin/bash
2  # $ -cwd
```

```
3  #$ -j y
4  #$ -l highmem
5  #$ -pe smp 6
6  #$ -l h_vmem=3G
7  #$ -l h_rt=12:0:0
8
9  module load anaconda3
10
11  conda activate mikado2.1
12
13  mikado prepare --json-conf configuration.yaml
```

The output of this script is two files: `mikado_prepared.fasta` and `mikado_prepared.gtf`.

### 3.3 Generate homology information and potential ORFs

Before generating a consensus set of genes/proteins, it is recommended to generate homology evidence for each of the predicted transcripts provided to Mikado, so that the program can use that information to improve annotation. To speed up the process, one can use [DIAMOND Blast](#). Apocrita has an old version, better install the newest stable via `conda`:

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -l highmem
5  #$ -pe smp 10
6  #$ -l h_vmem=4G
7  #$ -l h_rt=1:0:0
8
9  module load anaconda3
10  conda activate diamond
```

```
11
12 diamond blastx --query ../mikado_prepared.fasta --max-target-seq
    s 5 --sensitive --index-chunks 1 --threads 10 --db /data/SBCS-Ma
    rtinDuranLab/00-BlastDBs/Uniprot_SwissProt_160719.dmnd --evaluate
    1e-6 --outfmt 5 --out mikado.diamond.xml
```

The output is `mikado.diamond.xml`, which we can use in the final steps of Mikado.

In parallel, generate ORF predictions of the `mikado_prepared.fasta` transcripts employing [Transdecoder](#) (installed in my `$PATH`):

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -l highmem
5  #$ -pe smp 4
6  #$ -l h_vmem=3G
7  #$ -l h_rt=8:0:0
8
9  # Load the application module
10 module load perl
11
12 TransDecoder.LongOrfs -t mikado_prepared.fasta
13
14 TransDecoder.Predict -t mikado_prepared.fasta
```

Transdecoder generates a series of files, among them `mikado_prepared.fasta.transdecoder.bed`, which contains the predicted ORFs from the (redundant) genes initially generated by Mikado.

### 3.4 Merge all info and generate final output consensus gene set

As a final step, Mikado will take all evidences (input alignments, BLAST homologues,

ORF predictions, etc) to generate best-fit gene models. To do so, one needs to run the Mikado steps `serialise` and `pick`:

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -pe smp 1
5  #$ -l highmem
6  #$ -l h_vmem=10G
7  #$ -l h_rt=2:0:0
8
9  # Activate the genome annotation pipeline environment installed
   by James and add all modules into the $PATH
10 module load anaconda3
11 conda activate mikado2.1
12
13 # serialise and pick
14 mikado serialise --procs 1 --json-conf configuration.yaml --xml
   ./00-Diamond/mikado.diamond.xml --orfs ./01-Transdecoder/mikado_
   prepared.fasta.transdecoder.bed --blast_targets /data/SBCS-Marti
   nDuranLab/00-BlastDBs/Uniprot_SwissProt_160719.fasta --transcrip
   ts mikado_prepared.fasta --junctions ../00-Mikado_input/Owenia_j
   unctions_consensus.bed
15
16 mikado pick --procs 40 --json-conf configuration.yaml
```

The final output of this step is `mikado.loci.gff3` file, with all the transfrags (~ genes, isoforms, non coding RNAs, super loci, etc) considered by Mikado. This is a curated RNA-seq dataset that can be used as hints for Augustus ([Step 4](#)) and to update the annotation with PASA ([Step 5](#)).

**Note1:** With the update of Diamond, Mikado now gives an error for the `.xml` format of Diamond. This is probably a bug, since Mikado seems to load queries, targets and

alignments properly, and the BLAST metrics and scores are incorporated during the pick stage.

**Note2:** My experience is that it is better to use just 1 processor for `mikado` `serialise`.

## Step 4: Generating gene predictions with Augustus

[Augustus](#) is a pipeline that can either generate *ab initio* gene predictions (without any sort of guidance) or incorporate curated gene/protein/intron hints to generate high confident gene models. We will take advantage of these functionalities to train and run Augustus using our Mikado gene set ([Step 3](#)) and Portcullis exon junctions ([Step 2](#)). In addition, we will provide spliced-aligned proteins from a related species. A complete pipeline of how to use Augustus can be found [here](#).

### 4.1 Training Augustus

#### 4.1.1 First round of training

Use the script `select_mik_train.py` ([link](#)) to extract (i) full-length (start and stop codon) (ii) non-redundant transcripts with (iii) a blast hit (on a given fraction of their length, flag `-f`) and (iv) at least *n* exons (at least 2, flag `-e`):

```
select_mik_train.py -f 0.5 -e 2 mikado.loci.metrics.tsv mikado.l  
oci.gff3
```

This will generate a file called `training.gff3` containing the coding transcripts selected for training Augustus.

For training, Augustus needs to have full access to the `/path/to/Augustus/config` directory. Since the HPC installation does not give local access, I recommend downloading Augustus to your local `src` folder (which should be in your `PATH`) and pointing Augustus to that folder, but still running the HPC installation.

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

```
4  #$ -pe smp 4
5  #$ -l h_vmem=1G
6  #$ -l h_rt=48:0:0
7
8  # Load the application module
9  module load augustus/3.2.3
10
11 # Copy genome into working directory
12 cp /data/SBCS-MartinDuranLab/03-OweniaGenome/00-DATA/Owenia_soft
   masked_v072019.fa ./
13
14 # Set environmental variables and run training employing an Augu
   stus script that does the entire training pipeline
15 export AUGUSTUS_CONFIG_PATH=/data/SBCS-MartinDuranLab/02-Chema/s
   rc/Augustus/config
16 mkdir aug_training
17 autoAugTrain.pl --trainingset=training.gff3 --genome=Owenia_sof
   tmasked_v072019.fa --species=Owenia --workingdir=aug_training --
   optrounds=1 --verbose
```

For the first time, it is better to apply one round of optimisation and check the sensitivity and specificity at the nucleotide and exon level at the end of the training report, which is in `aug_training/autoAugTrain/training/test/augustus.2.withoutCRF.out`. If these parameters are not satisfactory, you can consider going for 2 rounds of optimisation and using CRF training instead of HMM (but check discussion regarding CRF training [here](#)).

The results after one round of optimisation are at the nucleotide level 94.5% of sensitivity and 26.9% of specificity, and at the exon level 83.2% of sensitivity and 38.3% of specificity. However, gene level prediction is very low: 18.5% sensitivity and 4.28% specificity. This might have something to do with the STOP codon, which might not be part of the last exon (see [here](#) to an explanation)

## 4.2 Generate gene hints

### 4.2.1 Generate species-specific exon hints

We can use the script `gtfToHintsMik.py` ([link](#)) to convert `mikado.loci.gff3` into hints.

```
1 cp /data/SBCS-MartinDuranLab/03-OweniaGenome/10-Mikado/mikado.loci.gff3 ./
2 # Convert gff3 to gtf (-T flag) and keep only coding transfrag
  (-C flag), highlighting any error (-E flag)
3 gffread -T -E -C -o mikado.loci.gtf mikado.loci.gff3
4 ./gtfToHintsMik.py mikado.loci.gtf
```

The output is `mikado.loci.exh.gff` (exh = exon hints)

### 4.2.2 Generate species-specific intron hints

We can use the information generated by Portcullis, to filter the original BAM file with all the RNA-seq data and generate a filtered BAM with reliable junctions. Then, we call an Augustus script to generate intron hints from the filtered BAM file of each library.

```
1 #!/bin/bash
2 #$ -cwd
3 #$ -j y
4 #$ -pe smp 4
5 #$ -l highmem
6 #$ -l h_vmem=1G
7 #$ -l h_rt=6:0:0
8 #$ -t 1-34
9
10 # Load the application module
11 module load augustus/3.2.3
12
13 # Get intron hints
14 bam2hints --intronsonly --minintronlen=15 --in=/data/scratch/btx
```



```
333/01-Annotation/08-STAR/03-Portcullis/portcullis_library_${SGE_TASK_ID}/portcullis.filtered.bam --out=/data/scratch/btx333/01-Annotation/11-Augustus/00-Hints/Owenia.intronthints_${SGE_TASK_ID}.gff
```

This generates 34 `Owenia.intronthints.gff`

#### 4.2.3 Generate spliced protein alignments with Exonerate

*Capitella* has a well-conserved gene repertoire, and as such, aligning its proteome (without TEs, filtered in [Step 1 after running RepeatModeler](#)) with [Exonerate](#) can provide evidences of genes not present in the RNA-seq data and/or support those obtained from species-specific RNA-seq/transcriptome alignments. Exonerate is slow. Before running it, split the fasta in multiple files with `fastasplit`

```
1 module load exonerate/2.4.0
2 fastasplit -f cte_prot_NR_noTE.fasta -o ./ -c 100
```

Now we can create an array job that takes it chunk and runs Exonerate on it. This involves renaming the chunks to something like this

`cte_prot_NR_noTE.fasta_chunk_1` with the last number being the number of the chunk, from 1 to 100. That will help build the array of jobs. Dividing the proteome in 100 chunks generates files of <400 proteins, which is quite efficient.

```
1 #!/bin/bash
2 # $ -cwd
3 # $ -pe smp 4
4 # $ -l h_vmem=4G
5 # $ -j y
6 # $ -l h_rt=6:00:00
7 # $ -t 1-100
8
9 # Load the application module
10 module load exonerate/2.4.0
11
```

```

12 # Run exonerate
13 exonerate --model protein2genome --cores 4 --bestn 3 --showtarget
    tgff T --showvulgar F --maxintron 500000 --fsmmemory 16000 --soft
    masktarget T cte_prot_NR_noTE.fasta_chunk_${SGE_TASK_ID} Owenia
    _softmasked_v072019.fa > Owenia_CapitellaProt_Exo_${SGE_TASK_I
    D}.gff

```

Merge all `.gff` files generated by Exonerate into one and extract a hint file employing Ferdi's script `exoToHints.py`.

```

1 cat Owenia_CapitellaProt_Exo_*.gff > Owenia_CapitellaProt_Exo.gff
2 grep "^Scaffold" Owenia_CapitellaProt_Exo.gff > Owenia_Capitella
  Prot_Exo_parsed.gff
3 ./exoToHints.py Owenia_CapitellaProt_Exo_parsed.gff

```

This generates `Owenia_CapitellaProt_Exo_parsed.exh.gff` (exh = exonerate hints)

#### 4.2.4 Merge all hints together and build a configuration file

All the steps to generate hints include in the final outputs the type of hint and the source, meaning that one can concatenate all of them in one single file.

```

cat mikado.loci.exh.gff Owenia_CapitellaProt_Exo_parsed.exh.gff
Owenia.intronthints.gff > Owenia_merged_hints.gff

```

After merging, one needs to modify the configuration file accordingly:

```

1 [GENERAL]
2   exonpart      100   .7   E 1  1e+3   P 1      1
3   intron        100   .7   E 1  1e+3   P 1      1
4   CDSpart       10    1    E 1      1    P 1     10

```

### 4.3 Augustus run

## Run Augustus:

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -l highmem
5  #$ -pe smp 30
6  #$ -l h_vmem=3G
7  #$ -l h_rt=72:0:0
8
9  # Load the application module
10 module load augustus/3.2.3
11
12 # Run augustus
13 export AUGUSTUS_CONFIG_PATH=/data/SBCS-MartinDuranLab/02-Chema/s
rc/Augustus/config
14
15 augustus --uniqueGeneId=true --gff3=on --species=Owenia --hintsf
ile=Owenia_merged_hints.gff --extrinsicCfgFile=extrinsic.Ferdi.
E.W.P.cfg --allow_hinted_splicesites=atac --alternatives-from-ev
idence=false Owenia_softmasked_v082020.fa > Owenia.aug.out
```

Augustus predicted 28,141 protein coding genes.

## Step 5: Merging Augustus and Mikado into a single gene set

To do so, we load the Augustus predicted proteins into [PASA](#) (which are based on curated hints and are solid predictions), and thereafter update them with the Mikado output, which is based on RNA-seq and thus contains UTR information, isoforms, etc.

Follow the instructions in the website to install PASA. In order to install the required

Perl modules locally, use [cpanm](#) and [local::lib](#). To use PASA with a SQLite database (easier to install and manage in Apocrita), you need to use the configuration file specifying the entire path to the database to be created. For instance, a standard alignment configuration file would look like this:

```
1  ## templated variables to be replaced exist as <__var_name__>
2
3  # database settings
4  DATABASE=/data/scratch/btx333/12-PASA/OweniaPASA_Augustus1st
5
6  #####
7  # Parameters to specify to specific scripts in pipeline
8  # create a key = "script_name" + ":" + "parameter"
9  # assign a value as done above.
10
11 #script validate_alignments_in_db.dbi
12 validate_alignments_in_db.dbi:--MIN_PERCENT_ALIGNED=<__MIN_PERCE
13 NT_ALIGNED__>
14
15 validate_alignments_in_db.dbi:--MIN_AVG_PER_ID=<__MIN_AVG_PER_ID
16 __>
17
18 #script subcluster_builder.dbi
19 subcluster_builder.dbi:-m=50
```

## 5.1 Extract protein coding genes from Augustus loci

One can extract the Augustus-based protein coding transcripts (CDS) using `gffread`:

```
gffread Owenia.aug.out -V -w Owenia.Augustus.genes.fa -g Owenia_
softmasked_v072019.fa
```

## 5.2 Run SeqClean on the Augustus coding genes

First, run `seqclean`. There is no need to include vector sequences, because the transcripts do not derive from a *de novo* assembly. `seqclean` will treat polyA regions, if included in the Augustus sets, and that will aid to better annotate terminating sites.

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -l highmem
5  #$ -pe smp 4
6  #$ -l h_vmem=2G
7  #$ -l h_rt=4:0:0
8
9  module load perl
10 module load samtools/1.9
11
12 /data/SBCS-MartinDuranLab/02-Chema/src/PASApipeline-v2.3.3/bin/s
   eqclean Owenia.Augustus.genes.fa
```

The output is `Owenia.Augustus.genes.fa.clean`

## 5.3 Run PASA pipeline with the clean transcripts

We call the main `PASA` pipeline to generate the initial database with the transcripts, including the following flags: `-C` creates the MYSQL database, `-R` runs the alignment/assembly pipeline, `-g` inputs the genome, `-t` inputs the cleaned transcripts, `-c` inputs the configuration file, `-u` the original transcripts and `-T` to indicate that the transcripts were trimmed using the TGI `seqclean` tool. You can use a standard `alignAssembly.config` file (as [above](#))

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

```
5  #$ -pe smp 40
6  #$ -l h_vmem=1G
7  #$ -l h_rt=18:0:0
8
9  module load perl
10 module load samtools/1.9
11
12 /data/SBCS-MartinDuranLab/02-Chema/src/PASApipeline-v2.3.3/Launc
   h_PASA_pipeline.pl -c alignAssembly.config -C -R -g Owenia_softm
   asked_v082020.fa -t Owenia.Augustus.genes.fa.clean -T -u Owenia.
   Augustus.genes.fa --ALIGNERS blat --CPU 40
```

## 5.4 Load the Mikado predictions to the PASA database

Once the database is created, we can incorporate the loci of Mikado. However, Mikado predicts both coding and non-coding genes, and the latter cannot be loaded into PASA. First, we need to “clean” the Mikado GFF3 file, to remove some non-canonical terms (e.g. super loci). We can do that with AGAT pipeline:

```
1  module load anaconda3
2  conda activate AGAT
3  agat_convert_sp_gxf2gxf.pl -g mikado.loci.gff3 -o mikado.loci.AG
   AT.gff3
```

Then, we need to remove the non-protein coding genes, which appear as ncRNA\_gene feature. We can grep the IDs of this genes, and then remove them from the GFF3 file with fgrep (much much faster than grep when working with large files).

```
fgrep -v -w -f mikado.ncRNA.IDs mikado.loci.AGAT.gff3 > mikado.l
oci.AGAT.NOncRNA.gff3
```

Once this is done, we can load the Mikado coding loci:

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -pe smp 2
5  #$ -l highmem
6  #$ -l h_vmem=1G
7  #$ -l h_rt=18:0:0
8
9  module load perl
10 module load samtools/1.9
11
12 /data/SBCS-MartinDuranLab/02-Chema/src/PASApipeline-v2.3.3/scripts/Load_Current_Gene_Annotations.dbi -c alignAssembly.config -g
Owenia_softmasked_v082020.fa -P mikado.loci.AGAT.NOncRNA.gff3
```

## 5.5 Update the PASA database with Mikado prediction

Next, compare the two annotations (the initially loaded based on Mikado and the *ab initio* of Augustus) and generate an updated, consensus one.

```
1  #!/bin/bash
2  #$ -cwd
3  #$ -j y
4  #$ -l highmem
5  #$ -pe smp 30
6  #$ -l h_vmem=1G
7  #$ -l h_rt=36:0:0
8
9  module load perl
10 module load samtools/1.9
11
12 /data/SBCS-MartinDuranLab/02-Chema/src/PASApipeline-v2.3.3/Launc
```

```
h_PASA_pipeline.pl -c annotCompare.config -A -g Owenia_softmaske  
d_v082020.fa -t Owenia.Augustus.genes.fa.clean --CPU 30
```

The output will have the format . Repeat the steps [5.4](#) and [5.5](#) again, starting with the output of the first update, to incorporate and update as many possible genes (recommended by PASA). The final output of these two rounds of updating the annotation is

OweniaPASA\_Augustus1st.gene\_structures\_post\_PASA\_updates.5127.gff3. We can use `AGAT` to calculate some statistics:

Type	Number
Gene	35,279
Transcript	40,268
5' UTR	36,222
Exon	265,085
CDS	250,227
3' UTR	33,349

## Step 6: Filtering Gene Models

Now that we have a merged all gene annotations, and before we start analysing this gene repertoire, we need to remove (i) spurious gene models; and (ii) genes that correspond to repeats and transposable elements (TEs) (annotated in [Step1](#)).

### 6.1 Verify and remove predictions with in-frame STOP codons

We can use again `AGAT` to clean that `GFF3` file, removing identical isoforms (it finds 4):

```
agat_convert_sp_gxf2gxf.pl -g OweniaPASA_Augustus1st.gene_struct  
ures_post_PASA_updates.5127.gff3 -o OweniaPASA.AGAT.gff3
```

The `gffread` can be used to identify any potential gene model with in frame STOP



codons (`-v` flag and `-H` to try to correct them):

```
gffread -E OweniaPASA.AGAT.gff3 -g Owenia_softmasked_v082020.fa
-V -H -o Owenia.PASA.simp.noSTOP.gff3
```

In this case, `gffread` finds 1 gene model (`mikado.Scaffold06G1968`) with in frame STOP codon. I remove it by hand and generate the file `OweniaPASA.AGAT.noSTOP.gff3`.

## 6.2 Remove gene predictions that overlap with repeats

First, we need to convert the repeat masker output file to a BED file, using Ferdi's [rep2bed.py](#) script:

```
1 module load python/2.7.15
2 ./rep2bed.py owenia_unmasked.fa.out > Ofus_RM.bed
```

Then, we can use Ferdi's [filt-rep-gtf.py](#) script, which requires the dependency `pybedtools` (easily installed with conda) to remove those gene models that overlap with repeats. Ferdi's script only works with `gtf` file formats, thus first thing is to convert the `gff3` to `gtf` (e.g. with `gffread`)

```
1 module load python/2.7.15
2 module load bedtools
3 pip install pybedtools --user #NOTE! This you only need to run o
nce, on the first time
4 ./filt-rep-gtf.py OweniaPASA_Augustus1st.AGAT.noSTOP.2.1.gtf Ofu
s_RM.bed
5
6 ## 250183 exons...
7 ## initial number of genes: 35278
8 ## number of genes after filtering 29385
```

We remove 5,893 gene models that overlap significantly with masked regions. The script generates two outputs, a "clean" `GTF` file and a list of retained genes. You can

use AGAT to obtain basic statistics of the final output:

```
agat_sq_stat_basic.pl -i OweniaPASA.AGAT.noSTOP.filt.AGAT.gff3 -
g Owenia_softmasked_v082020.fa
```

Type	Number	Size mean (kb)	% genome
Gene	29,385	12,009.98	70.56
Transcript	34,353	13,547.80	-
5' UTR	39,026	137.88	1.08
Exon	257,455	278.37	14.33
CDS	242,753	223.50	10.85
3' UTR	32,941	365.24	2.41

## 6.3 Remove gene predictions that hit TEs

Finally, we could have predicted TEs expressed in our RNA-seq datasets. We can implement a similar strategy than the one followed to de-contaminate *Capitella's* proteome in [Step 1](#). First, we extract the protein coding genes and then we BLAST them against a library of TE and repetitive elements. Finally, we remove those gene models that give hit against TEs.

```
1 gffread -E OweniaPASA.AGAT.noSTOP.filt.gtf -S -g Owenia_softmask
   ed_v082020.fa -y OweniaPASA.AGAT.noSTOP.filt.prot.fasta
2 cp /data/SBCS-MartinDuranLab/02-Chema/00-OweniaGenome/02-Filteri
   ngRB/RepeatPeps.dmnd ./
3 conda activate diamond
4 diamond makedb --in RepeatPeps.lib -d RepeatPeps
5 diamond blastp -d RepeatPeps -q OweniaPASA.AGAT.noSTOP.filt.pro
   t.fasta -o Owenia.VS.TEs.1e5.blastp -f 6 qseqid bitscore evalue
   stitle -k 25 -e 1e-5 -p 8
6 ## 2450 queries aligned.
```

We identify 2,450 gene models as potential TEs. Most of them have very low e-values, meaning that are true TEs. Those with higher e-values are small predictions. I manually checked a couple of them, and they give hit against TEs. To remove them, we first get the gene ID and then we `fgrep` them out (with the `-v` flag):

```
1 cat Owenia.VS.TEs.1e5.blastp | cut -f 1 | sort | uniq > TEsIDs
2 fgrep -w -v -f TEsIDs OweniaPASA.AGAT.noSTOP.filt.gtf > OweniaPA
  SA.AGAT.noSTOP.filt.noTEs.gtf
```

I re-run the analysis to make sure that all TEs are gone:

```
1 gffread -E Owenia.PASA.simp.noSTOP.noRepeat.noTE.gtf -S -g Oweni
  a_softmasked_v072019.fa -y Owenia.PASA.simp.noSTOP.noRepeat.noT
  e.prot.fasta
2 diamond blastp -d RepeatPeps -q Owenia.PASA.simp.noSTOP.noRepea
  t.noTe.prot.fasta -o OweniaPASAnoSTOPnoRepeatnoTE.vs.RepeatPeps.
  1e5.blastp -f 6 qseqid bitscore evalue stitle -k 25 -e 1e-5 -p 8
3 ## 0 queries aligned.
```

Because the GTF format (in version 2) does not retain UTR information etc, we can add that information with AGAT, transforming it to `GFF3`:

```
agat_convert_sp_gxf2gxf.pl -g OweniaPASA.AGAT.noSTOP.filt.gtf -o
  OweniaPASA.AGAT.noSTOP.filt.AGAT.gff3
```

You can use AGAT to obtain basic statistics of the final output:

```
agat_sq_stat_basic.pl -i OweniaPASA.AGAT.noSTOP.filt.noTE.AGAT.g
  ff3 -g Owenia_softmasked_v082020.fa
```

Type	Number	Size mean (kb)	% genome
Gene	26,966	12,763.72	68.82

Transcript	31,903	14,305.44	91.25
5' UTR	37,141	134.57	1.00
Exon	252,326	262.47	13.24
CDS	237,854	209.32	9.95
3' UTR	31,076	368.17	2.29

## 6.4 Rename annotation file, gene predictions and produce a Non-Redundant version

At this point, save the original PASA file

(OweniaPASA\_Augustus1st.gene\_structures\_post\_PASA\_updates.5127.gff3) and the filtered annotation () to the lab folder.

For the sake of simplicity with downstream analyses, we can rename gene names. We can do that with AGAT, adopting an Ensembl like nomenclature:

```
agat_sp_manage_IDs.pl -f OweniaPASA_Augustus1st.AGAT.noSTOP.noRe
peats.gff3 --ensembl --prefix OFUS --type_dependent --tair -o Ow
enia_annotation_v250920.gff3
```

This command generates a final output named Owenia\_annotation\_v250920.gff3 that will be used in all downstream analyses. We can additionally generate an annotation file without isoforms (picking just the longest one), with AGAT, which will be used in gene family evolution analyses (see [+Gene family evolution: Step-2:-Generate-Non-Redundant](#))

```
1 agat_convert_sp_gxf2gxf.pl --gff Owenia_annotation_v250920.gff3
  --merge_loci -o Owenia_lociMerged.gff
2 agat_sp_keep_longest_isoform.pl --gff Owenia_lociMerged.gff -o O
  wenia_lociMerged_longestIsoform.gff
```

The Ensembl nomenclature has potential for 11 digits numbers, so I just adjust that with a text editor so that it is 5-digit based (i.e. OFUSG00001 and OFUSG13589). The transcript/isoform is indicated after the dot, i.e. OFUSG00001.1 and

OFUSG00001.2 are isoforms 1 and 2 of the gene OFUSG00001 respectively.

## Step 7: Validating the annotation

### 7.1 Run BUSCO

We can use BUSCO to assess that the final annotated gene set has similar completeness than the entire genome (see [Step 0](#)).

```
gffread -E Owenia_annotation_v250920.gff3 -g Owenia_softmasked_v082020.fa -y Owenia_annotation_v250920_filteredPeps.fasta
```

```
1  #!/bin/bash
2  #$ -pe smp 20
3  #$ -l highmem
4  #$ -l h_vmem=1G
5  #$ -l h_rt=0:30:0
6  #$ -cwd
7  #$ -j y
8
9  module load anaconda3
10
11  conda activate BUSCO
12
13  busco -m proteins -c 20 -i Owenia_annotation_v250920_filteredPeps.fasta -o Owenia_final -l metazoa_odb10
14
15  conda deactivate
```

The results are very comparable to what was reported for the BUSCO genome (97.8% vs 99%)

```
1  -----
```

```

2      |Results from dataset metazoa_odb10      |
3      -----
4      |C:96.2%[S:86.9%,D:9.3%],F:1.5%,M:2.3%,n:954|
5      |918      Complete BUSCOs (C)           |
6      |829      Complete and single-copy BUSCOs (S)|
7      |89       Complete and duplicated BUSCOs (D)|
8      |14       Fragmented BUSCOs (F)          |
9      |22       Missing BUSCOs (M)            |
10     |954      Total BUSCO groups searched    |
11     -----

```

## Step 8: Generate a functionally annotated database

The final step of the annotation is to functionally annotate the gene set, i.e. assign gene orthology, PFAM domains, transmembrane regions, GO terms, etc to each gene model. To do so, we use [Trinotate](#).

### 8.1 Set up Trinotate

To set up Trinotate, you need to have two Perl5 modules installed locally `DBI` and `DBD::SQLite`. To install a perl module locally, you need to add a `PREFIX` flag when running `Makefile.PL`:

```

1  module load perl
2  cd <DIRECTORY WHERE THE PERL MODULE IS>
3  perl Makefile.PL PREFIX=~/.lib/perl5
4  make
5  make test
6  make install

```

Then, add to your `.bash_profile` the following line:

```
export PERL5LIB=~/lib/perl5/lib/site_perl
```

Follow guidelines in the [Trinotate wiki](#). Trinotate still uses some old versions (e.g. of signalP) and it is important to employ the proper versions of the different modules.

**Note:** I don't install RNAMMER, as it requires a lot of tinkering. Just install signalP and tmhmm in your user, and use the NCBI BLAST+ and HMMER already installed in Apocrita.

```
1 #Transcript ID
2 >Ofus.G0115.1
3 #Peptide ID
4 >Ofus.G0115.1 len:142 Ofus.G0115.1:1-429(+)
5 #GeneTransMap* **
6 Ofus.G0115.1 len:142 Ofus.G0115.1:1-429(+)          Ofus.G0115.1
7 peptideID<TAB>transcriptID
```

\*When the output comes from a Trinity assamble and not from the genome, in order to get this file one needs to use a line from the Trinity utilities directory

```
$TRINITY_HOME/util/support_scripts/get_Trinity_gene_to_trans_map.pl Trinity.fasta > Trinity.fasta.gene_trans_map
```

Apocrita specifics:

```
1 #!/bin/bash
2 #$ -pe smp 1
3 #$ -l highmem
4 #$ -l h_vmem=5G
5 #$ -l h_rt=2:0:0
6 #$ -cwd
7 #$ -j y
8
9 module load trinity/2.4.0
10 module load intel/2017.1
```

```

11 module load bowtie2/2.3.2
12
13 /data/apps/trinity/2.2.0/trinityrnaseq-2.2.0/util/support_scripts/
  get_Trinity_gene_to_trans_map.pl UrechisRefTrans_cdhit.fasta >
  UrechisRefTrans_cdhit.fasta.gene_trans_map

```

**\*\*or another option would be to get just the transcripts that come from coding genes (i.e. those who will make proteins). Therefore it is necessary to use the [#Peptide](#) ID file created from the transdecoder.pep above and obtain the transcript ID using grep and trimmed it to obtain something like this:**

```

1 >TRINITY_DN100003_c0_g1_i1.
2 >TRINITY_DN100006_c0_g1_i1.
3 >TRINITY_DN100008_c0_g1_i1.
4 >TRINITY_DN100015_c0_g1_i1.
5 etc...

```

and then use FastaGREG to obtain the sequences of only the coding genes from the Trinity assembly. Example below is with the RNA-seq from Urechis (Park et al.).  
**SUBMIT as a Job**

```

1 #!/bin/bash
2 #$ -cwd
3 #$ -j y
4 #$ -pe smp 8
5 #$ -l h_vmem=1G
6 #$ -l h_rt=36:0:0
7
8 module load perl
9 FastaGREG -f transcriptID -X UrechisRefTrans_cdhit.fasta > Urech
  isRefTrans_coding.fasta

```

This new fasta file will be the input for the option in Trinotate: `--transcript_fasta`

## 8.1 Giacomo

Installation with anaconda3:



```
1 module load anaconda3
2 conda create --prefix /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/trinotate_env
3 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/trinotate_env
4 conda install -c bioconda trinotate
```

First thing we need to do is to generate 2 fasta files from the final annotation gff3 file:

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/$species/annotation/final/$final_annotation
13 species_softmasked="$species"_softmasked.fa
14 species_softmasked_path=/data/SBCS-MartinDuranLab/03-Giacomo/data/$species/annotation/softmasking/$species_softmasked
15 output_mRNA="$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/
25
26 mkdir -p annotation_step8
27 cd annotation_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 -y $output_proteins -g $species_softmasked $final_annotation
```

Then we need to modify these two fasta files in order to make them recognisable by Trinotate:

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="$species".gene_trans_map
15
16 cd /data/scratch/btx654/btx604-scratch/$species/annotation_step8
17 /
18 grep ">" $fasta_mRNA | awk '{ print $2"(" + ")"}' > positions.txt
19 sed -i 's/CDS=/:/' positions.txt
20
21 awk '/^>/ {if (seqlen){print "len:"seqlen}; print ;seqlen=0;nex
t; } { seqlen += length($0)}END{print "len:"seqlen}' $fasta_prot
eins | grep -v ">" > lenghts.txt
22
23 grep ">" $fasta_proteins > names.txt
24 sed 's/> //' names.txt > names_clean.txt
25
26 sed '/^>/s/ .*//' $fasta_mRNA > input_trinotate_mRNA.fa
27
28 cp $fasta_proteins input_trinotate_proteins.fa
29 INDEX=1
30 while read -r line
31 do
32 original_name=$(echo $line)
33 lenght=$(head -$INDEX lenghts.txt | tail -1)
34 name_clean=$(head -$INDEX names_clean.txt | tail -1)
35 position=$(head -$INDEX positions.txt | tail -1)
36
37 sed -i "s/$original_name/$original_name $lenght $name_clean$posi
tion/" input_trinotate_proteins.fa
38
39 INDEX=$((INDEX+1))
40 done < names.txt
```

```
41
42
43 grep ">" input_trinotate_proteins.fa > full_names.txt
44 sed -i 's/> //' full_names.txt
45
46 INDEX=1
47 while read -r line
48 do
49     full_name=$(echo $line)
50     name_clean=$(head -${INDEX} names_clean.txt | tail -1)
51
52     echo -e $full_name'\t'$name_clean >> $gene_trans_map
53
54     INDEX=$((INDEX+1))
55 done < full_names.txt
```

## 8.2 Generate input files

Trinotate will consider five functional evidences:

- **BLASTp**

```
1 #!/bin/bash
2 # $ -cwd
3 # $ -j y
4 # $ -pe smp 8
5 # $ -l h_vmem=1G
6 # $ -l h_rt=36:0:0
7
8 module load blast+
9
10 blastp -query Owenia_filtered_proteins.fa -db /data/SBCS-MartinDuranLab/02-Chema/src/Trinotate-Trinotate-v3.2.0/TrinotateDBs/uni
```

```
prot_sprot.pep -num_threads 8 -max_target_seqs 1 -outfmt 6 -eval  
ue 1e-3 > blastp.outfmt6
```

### • BLASTx

```
1 #!/bin/bash  
2 #$ -cwd  
3 #$ -j y  
4 #$ -pe smp 8  
5 #$ -l h_vmem=1G  
6 #$ -l h_rt=36:0:0  
7  
8 module load blast+  
9  
10 blastx -query Owenia_filtered_CDS.fa -db /data/SBCS-MartinDuranL  
ab/02-Chema/src/Trinotate-Trinotate-v3.2.0/TrinotateDBs/uniprot_  
sprot.pep -num_threads 8 -max_target_seqs 1 -outfmt 6 -evaluate 1e  
-3 > blastx.outfmt6
```

### • HMMER

```
1 #!/bin/bash  
2 #$ -cwd  
3 #$ -j y  
4 #$ -l highmem  
5 #$ -pe smp 12  
6 #$ -l h_vmem=1G  
7 #$ -l h_rt=36:0:0  
8  
9 module load hmmer/  
10  
11 hmmscan --cpu 12 --domtblout PFAM.out /data/SBCS-MartinDuranLab/  
00-BlastDBs/Pfam-A.hmm ./00-DATA/Ofus_peps.fasta > pfam.log
```

### • SignalP

```
1 #!/bin/bash
2 #$ -cwd
3 #$ -j y
4 #$ -pe smp 1
5 #$ -l h_vmem=8G
6 #$ -l h_rt=36:0:0
7
8 module load perl
9
10 signalp -f short -n signalp.out Owenia_filtered_proteins.fa
```

## 8.2 Giacomo

Before running BLASTp, BLASTx, HMMER and SignalP we need to download the databases we will use:

download\_prepare\_trinotate\_databases\_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=36:0:0
8 #$ -l highmem
9
10 module load anaconda3
11 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
12 3/trinotate_env
13
14 cd /data/SBCS-MartinDuranLab/03-Giacomo/db/trinotate
15
16 /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda3/trinotate_env
```

```
16 /bin/Build_Trinotate_Boilerplate_SQLite_db.pl Trinotate
17 makeblastdb -in uniprot_sprot.pep -dbtype prot
18
19 gunzip Pfam-A.hmm.gz
20 hmmcompress Pfam-A.hmm
```

Then we need to install SignalP in our src folder (outside the conda environment).  
Download signalp from [here](#) and then scp it into apocrita

```
tar -zxvf signalp-4.1g.Linux.tar.gz
```

and then go to the new folder and edit with nano the executable "signalp". You need to specify the path to the signalp folder in one of the first lines:

```
$ENV{SIGNALP} = '/data/SBCS-MartinDuranLab/03-Giacomo/src/signalp-4.1';
```

in the same lines of the script increase max allowed entries to a big number

```
my $MAX_ALLOWED_ENTRIES=2000000;
```

Now we are ready to generate the functional evidences. In order to make this step even easier I wrote a wrapper script that will launch automatically the single softwares on a specified species (\$1). Remember to specify the path to the databases installed with download\_prepare\_trinotate\_databases\_v1.sh in each of the 4 scripts contained in the wrapper.

The following script should be submitted in this way:

```
qsub wrapper_trinotate_8.2_v1.sh riftia
```

wrapper\_trinotate\_8.2\_v1.sh

```
1 #!/bin/bash
2 #$ -wd /data/home/btx654/scripts/annotation/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 qsub BLASTp_universal.sh $species
11 qsub BLASTx_universal.sh $species
12 qsub HMMER_universal.sh $species
13 qsub Signalp_universal.sh $species
```

These scripts are the ones that are submitted by the previous script:

#### 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/annotation_step8
14 /
15
16 module load anaconda3
17
18 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
19 3/trinotate_env
20
21 blastp -query input_trinotate_proteins.fa -db /data/SBCS-MartinD
22 uranLab/03-Giacomo/db/trinotate/uniprot_sprot.pep -num_threads 8
23 -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/annotation_step8
14 /
15 module load anaconda3
16 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
17 3/trinotate_env
18 blastx -query input_trinotate_mRNA.fa -db /data/SBCS-MartinDuran
19 Lab/03-Giacomo/db/trinotate/uniprot_sprot.pep -num_threads 8 -ma
20 x_target_seqs 1 -outfmt 6 -evalue 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=1G
8  #$ -l h_rt=36:0:0
9
10 species=$1
11
12 echo "Working on "$species
13
14 cd /data/scratch/btx654/btx604-scratch/$species/annotation_step8
15 /
16
17 module load anaconda3
18
19 conda activate /data/SBCS-MartinDuranLab/03-Giacomo/src/anaconda
20 3/trinotate_env
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
```

### 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/annotation_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
```

### 8.3 Populate the SQLite database and output the report

```
1 /data/SBCS-MartinDuranLab/02-Chema/src/Trinotate-Trinotate-v3.2.
0/admin/Build_Trinotate_Boilerplate_SQLite_db.pl Owenia
2 /data/SBCS-MartinDuranLab/02-Chema/src/Trinotate-Trinotate-v3.2.
1/Trinotate Owenia.sqlite init --gene_trans_map 00-DATA/Ofus_gen
eTransMap --transcript_fasta 00-DATA/Ofus_mRNA.fasta --transdeco
der_pep 00-DATA/Ofus_peps.fasta
3 #LOAD annotations, below are examples from Trinotate manual, cha
nge accordingly to your species output from blast+:
4 /data/SBCS-MartinDuranLab/02-Chema/src/Trinotate-Trinotate-v3.2.
1/Trinotate Owenia.sqlite LOAD_swissprot_blastp blastp.outfmt6
5 /data/SBCS-MartinDuranLab/02-Chema/src/Trinotate-Trinotate-v3.2.
1/Trinotate Owenia.sqlite LOAD_swissprot_blastx blastx.outfmt6
6 /data/SBCS-MartinDuranLab/02-Chema/src/Trinotate-Trinotate-v3.2.
1/Trinotate Owenia.sqlite LOAD_signalp signalp.out
7 /data/SBCS-MartinDuranLab/02-Chema/src/Trinotate-Trinotate-v3.2.
1/Trinotate Owenia.sqlite LOAD_pfam TrinotatePFAM.out
8
9 /data/SBCS-MartinDuranLab/02-Chema/src/Trinotate-Trinotate-v3.2.
1/Trinotate Owenia.sqlite report --incl_pep --incl_trans > Oweni
a_annotation_report.xls
```

## Step 9: Generate PANTHER annotation

**Panther** is a classification system for proteins that combines the specificity of HMMR searches with a curated GO annotation for hundreds of genomes. It allow to classify genes in 15,702 protein families, divided into 123,989 functionally distinct protein subfamilies, all of them with associated GO terms. In order to run PANTHER and annotate each gene of the genome, one can download the Panther HMM Scoring tool from its [website](#) and run it like so (add the `-s` flag so that it uses HMMscan instead of HMMsearch, otherwise it will never finish). Be sure to download the script itself "pantherScore2.2.pl" and the folder "lib" containing the required perl modules. The directory "lib" has to be in the same folder of the script "pantherScore2.2.pl". Finally this script requires also the PANTHER15.0 directory downloadable from the same link (which it is needed for the `-l` option).

```
1  #!/bin/bash
2  #$ -pe smp 30
3  #$ -l highmem
4  #$ -l h_vmem=1G
5  #$ -l h_rt=120:00:0
6  #$ -cwd
7  #$ -j y
8
9  module load perl
10 module load hmmer/
11
12 pantherScore2.2.pl -l /data/SBCS-MartinDuranLab/00-BlastDBs/asci
i/PANTHER15.0/ -D B -n -o Owenia_Panther -i Owenia_annotation_v2
50920_filteredPeps.fasta -c 30 -V -s
```

**Note:** `pantherScore` is really really slow, give it >5 days to run for a genome of ~30,000 proteins!

Once it is done, you can combine Trinotate and Panther in a single output. To do so:

```
1  cut -f 1 Owenia_Panther_sorted > IDs_panther
2  cut -f 2 Owenia_annotation_v250920_report.xls | tail -n +2 > IDs
```

```

_all
3 fgrep -v -f IDs_panther IDs_all > IDs_absentPanther ### There are
  4923 genes without Panther annotation
4 awk '{print $0"\t""NO PTHR""\t""NO HIT"}' IDs_absentPanther > PAN
  THER_nohits
5 cat Owenia_Panther_sorted PANTHER_nohits | sort -k 1,1 > Owenia_
  Panther_sorted_allgenes
6 ## use vim to add a header in Owenia_Panther_sorted_allgenes so
  that it matches Trinotate file
7 paste Owenia_annotation_v250920_report.xls Owenia_Panther_sorted
  _allgenes > Owenia_annotation_v250920_TrinoPanther.xls

```

There are some minor formatting issues (e.g. there were a bunch of duplicated Panther terms that misaligned the correspondence between Trinotate and Panther annotations) that one can correct by hand easily directly on Excel.

## 9 Giacomo

### 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 16/11/20
12 gffread_proteins_path=/data/scratch/btx654/btx604-scratch/$speci
  es/annotation_step8/$gffread_proteins
13

```

```
14 echo "Working on "$species
15
16 cd /data/scratch/btx654/btx604-scratch/$species/
17
18 mkdir -p annotation_step9
19 cd annotation_step9
20
21 cp $gffread_proteins_path ./
22
23 module load perl
24 module load hmmer/
25
26 export PERL5LIB=/data/SBCS-MartinDuranLab/03-Giacomo/src/hmmscoring/lib/
27
28 perl /data/SBCS-MartinDuranLab/03-Giacomo/src/hmmscoring/panther
Score2.2.pl -l /data/SBCS-MartinDuranLab/03-Giacomo/src/hmmscoring/PANTHER15.0/ -D B -n -o riftia_panther -i $gffread_proteins -
c 30 -V -s
```

## Step 10. Lift-over (if needed)

At some point, one might need to transfer the annotation from one assembly version to another (better) one. One can do that with [Liftoff](#) pipeline like so:

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

```
7  # $ -j y
8
9  module load anaconda3
10
11  conda activate Liftoff
12
13  liftoff -p 5 -g Owenia_annotation_v300321.1_filtered_noGeneID.gff3 -o Owenia_chrom_v300321.gff3 Owenia_v082020_hic.fa Owenia_unmasked_v082020.fa
14
15  conda deactivate
```