

Pipelines of Marker Development for Transcriptome-based Exon Capture

Part I phylogenomics

January 17, 2015

Contributors: Sonal Singhal and Ke Bi

For questions or to report bugs, please contact Ke Bi (kebi@berkeley.edu)

Reference:

- [1]. Singhal S. 2013. De novo transcriptomic analyses for non-model organisms: an evaluation of methods across a multi-species data set. *Molecular Ecology Resources* 13:403-416.
- [2]. Bi K, Linderoth T, Vanderpool D, Good JM, Nielsen R and Moritz C. 2013. Unlocking the vault: next-generation museum population genomics. *Molecular Ecology* 22:6018-6032.
- [3]. Bi K, Vanderpool D, Singhal S, Linderoth T, Moritz C and Good JM. 2012. Transcriptome-based exon capture enables highly cost-effective comparative genomic data collection at moderate evolutionary scales. *BMC Genomics* 13: e403.

The pipelines are deposited in
<https://github.com/CGRL-QB3-UCBerkeley/MarkerDevelopmentPylogenomics>

Scripts included in this pipeline:

[1-PreCleanup](#)

[2-ScrubReads](#)

[3-GenerateAssemblies](#)

[4-AssemblyEvaluation](#)

[5-Annotation](#)

[6-MarkerSelectionTRANS](#)

[6-MarkerSelectionEXONS](#)

**Use "chmod +x script" to make each of these perl scripts executable.

46 **Note: If exon identification is not possible or not desirable, users can use the
48 entire transcripts for marker development. In this case please use "6-
MarkerSelectionTRANS". Otherwise please use "6- MarkerSelectionEXONS".

50

52 **1-PreCleanup**: Reformats raw cDNA sequencing reads from Illumina HiSeq or
MiSeq for [2-ScrubReads](#). Specifically, in this step we will remove reads that did not
54 pass the Illumina quality control filters and modify the sequence identifiers.

56 Dependencies:

FastQC: <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

58

Input:

60 Raw sequence data files are grouped and saved in folders named by their sample
IDs. For instance, three libraries (CGRL_index1, CGRL_index15, CGRL_index40) are
62 saved under “/home/ke/Desktop/SeqCap/data/rawdata/library/”. Compressed
fastq sequence files are saved in each of these folders.

64

Fastq files use the following naming scheme:

66 <sample name>_<barcode sequence>_L<lane (0-padded to 3 digits)>_R<read
number>_<set number (0-padded to 3 digits)>.fastq.gz

68

For example, in “CGRL_index15_CGACCTG_L006_R1_001.fastq.gz”:

70 sample name: CGRL_index15

barcode sequence: CGACCTG

72 lane (0-padded to 3 digits): 006

read number: 1

74 set number (0-padded to 3 digits): 001

76 #Make a new folder called “raw” under

“~/Desktop/MarkerDevelopment/data/rawdata/”:

78 *ke@NGS:~/Desktop/MarkerDevelopment/data/rawdata\$ mkdir raw*

80 #Copy all these compressed fastq files from each folder (CGRL_index1,
CGRL_index15, CGRL_index40) to “raw”:

82 *ke@NGS:~/Desktop/MarkerDevelopment/data/rawdata\$ cp*
library/CGRL_index/*.gz raw/*

84

#Check data files in “raw”:

86 *ke@NGS:~/Desktop/MarkerDevelopment/data/rawdata\$ ls raw/**

CGRL_index15_CGACCTG_L006_R1_001.fastq.gz

88 *CGRL_index15_CGACCTG_L006_R2_001.fastq.gz*

CGRL_index1_TCGCAGG_L006_R1_001.fastq.gz

90 *CGRL_index1_TCGCAGG_L006_R2_001.fastq.gz*

CGRL_index40_TTCGCAA_L006_R1_001.fastq.gz

92 *CGRL_index40_TTCGCAA_L006_R2_001.fastq.gz*

94

Commands:

```

96  #cd to the working directory:
    ke@NGS:~/Desktop/MarkerDevelopment/data/rawdata$ cd ..
98
    #run 1-PreCleanup with fastq evaluation
100 ke@NGS:~/Desktop/MarkerDevelopment/data$ 1-PreCleanup
    ~/Desktop/MarkerDevelopment/data/rawdata/raw/ fastqc
102
    ~/Desktop/MarkerDevelopment/data
104 Output:
    Three new folders will be created under
106 "~/Desktop/MarkerDevelopment/data/rawdata/raw/":
    "pre-clean"
108 "combined"
    "pre-clean/evaluation"
110
    - Folder "pre-clean" contains reformatted raw fastq reads.
112 CGRL_index1_R1.fq
    CGRL_index1_R2.fq
114 CGRL_index15_R1.fq
    CGRL_index15_R2.fq
116 CGRL_index40_R1.fq
    CGRL_index40_R2.fq
118
    - Folder "combined" contains merged, compressed, fastq data files (not used by the
120 following pipeline).
    CGRL_index1_TCGCAGG_L006_R1.fastq.gz
122 CGRL_index1_TCGCAGG_L006_R2.fastq.gz
    CGRL_index15_CGACCTG_L006_R1.fastq.gz
124 CGRL_index15_CGACCTG_L006_R2.fastq.gz
    CGRL_index40_TTCGCAA_L006_R1.fastq.gz
126 CGRL_index40_TTCGCAA_L006_R2.fastq.gz

128 - Folder "evaluation" contains fastQC results for each data file.
    CGRL_index1_R1.fq_fastqc/
130 CGRL_index1_R2.fq_fastqc/
    CGRL_index15_R1.fq_fastqc/
132 CGRL_index15_R2.fq_fastqc/
    CGRL_index40_R1.fq_fastqc/
134 CGRL_index40_R2.fq_fastqc/

```

```

136

```

138 *2-ScrubReads*: Clean up raw data, which includes trimming for quality, removing
 140 adapters, merging overlapping reads, removing duplicates and reads sourced from
 contamination

142 Dependencies:
 cutadapt: <http://code.google.com/p/cutadapt/>
 144 COPE: <http://sourceforge.net/projects/coperead/>
 Bowtie2: <http://sourceforge.net/projects/bowtie-bio/files/bowtie2/>
 146 FastQC: <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
 FLASh-modified: modified version of FLASh by Filipe G. Vieira.
 148 <https://github.com/MVZSEQ/Exon-capture>
 Trimmomatic: <http://www.usadellab.org/cms/?page=trimmomatic>

150

Input:

152 1. Reformatted fastq files created by [1-PreCleanup](#):
 #Check the raw data files:
 154 *ke@NGS:~/Desktop/MarkerDevelopment/data/rawdata/raw/pre-clean\$ ls *.fq*
CGRL_index1_R1.fq
 156 *CGRL_index1_R2.fq*
CGRL_index15_R1.fq
 158 *CGRL_index15_R2.fq*
CGRL_index40_R1.fq
 160 *CGRL_index40_R2.fq*

162 2. A fasta file that contains adapter sequences:
 #Check the format of adapter sequence file:
 164 *ke@NGS:~/Desktop/SeqCap/denovoTargetCapture/associated_files \$ less -S*
Adapters.fasta
 166 *>P7_index1*
CAAGCAGAAGACGGCATACGAGATcctgcgaGTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT
 168 *>P7_index2*
CAAGCAGAAGACGGCATACGAGATtgcagagGTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT
 170 *.....*
>P5_index1
 172 *AATGATACGGCGACCACCGAGATCTACACcctgcgaACACTCTTCCCTACACGACGCTCTTCCGATCT*
>P5_index2
 174 *AATGATACGGCGACCACCGAGATCTACACtgcagagACACTCTTCCCTACACGACGCTCTTCCGATCT*
 176 *.....*

Note: The header of each adapter sequence has to be named strictly as “**P7_indexN**”
 178 or “**P5_indexN**”. N is the number of index. It is OK to put all adapters in this file but
 your libraries only use a subset of them.

180

3. Library info file (Tab-delimited txt file):
 182 #Check the format of Library info file:
ke@NGS:~/Desktop/SeqCap/denovoTargetCapture/associated_files \$ less -S libInfo.txt
 184

<i>library</i>	<i>P7</i>	<i>P5</i>
----------------	-----------	-----------

186 *CGRL_index1* 1
 187 *CGRL_index15* 15
 188 *CGRL_index40* 40

190 Leave the “P5” column blank if you only have indexes in P7 adapters in the libraries.

192 4. Contaminant file:
 193 *Escherichia coli* (bacteria + human + other genome resources if desired) genome in
 194 fasta format.
 195 This file (e_coli_K12.fasta) is saved in
 196 “~/Desktop/SeqCap/denovoTargetCapture/associated_files/ecoli/”

198

Commands:

200 #Make a new folder called “cleaned_data” in
 201 “~/Desktop/MarkerDevelopment/data/”:
 202 ke@NGS:~/Desktop/MarkerDevelopment/data\$ *mkdir cleaned_data*

204 #Run [2-ScrubReads](#):
 205 ke@NGS:~/Desktop/MarkerDevelopment/data\$ *2-ScrubReads -f*
 206 *~/Desktop/MarkerDevelopment/data/rawdata/raw/pre-clean/ -o*
 207 *~/Desktop/MarkerDevelopment/data/cleaned_data/ -a*
 208 *~/Desktop/SeqCap/denovoTargetCapture/associated_files/Adapters.fasta -b*
 209 *~/Desktop/SeqCap/denovoTargetCapture/associated_files/libInfo.txt -t*
 210 */home/ke/Desktop/SeqCap/programs/Trimmomatic-0.32/trimmomatic-0.32.jar -c*
 211 *~/Desktop/SeqCap/denovoTargetCapture/associated_files/ecoli/e_coli_K12.fasta -e*
 212 *200 -m 15 -z*

214 Note: I use the default values for most of the arguments. Users should adjust these
 215 parameters when processing the real datasets.

216

Output:

218 1. In “~/Desktop/MarkerDevelopment/data/cleaned_data/”, six .txt files per
 219 library are produced:
 220 For example for library *CGRL_index1*, the six files are:
 221 *CGRL_index1_1_final.txt* (left reads)
 222 *CGRL_index1_2_final.txt* (right reads)
 223 *CGRL_index1_u_final.txt* (merged or unpaired reads)
 224 *CGRL_index1.contam.out* (headers of reads aligned to bacteria)
 225 *CGRL_index1.duplicates.out* (headers of duplicated reads)
 226 *CGRL_index1.lowComplexity.out* (headers of low complexity reads)

228 2. In “~/Desktop/MarkerDevelopment/data/cleaned_data/evaluation/”, you can
 229 find fastQC results for cleaned reads from each library.

230

3-GenerateAssemblies: Assemble RNAseq data using Trinity.

Dependencies:

Trinity <http://trinityrnaseq.sourceforge.net>

Input:

For each library, we will concatenate cleaned forward reads (XXX_1_final.txt) and unpaired reads (XXX_u_final.txt) and name the resulting read data file as XXX_1_final.txt.

#Make a new folder called "raw_assembly" under

"~/Desktop/MarkerDevelopment/data/":

```
ke@NGS:~/Desktop/MarkerDevelopment/data$ mkdir raw_assembly
```

#Concatenate cleaned forward reads and unpaired reads and save them in "raw_assembly":

```
ke@NGS:~/Desktop/MarkerDevelopment/data$ cat
```

```
cleaned_data/CGRL_index1_1_final.txt cleaned_data/CGRL_index1_u_final.txt | sed 's/\2$/\1/g' > raw_assembly/CGRL_index1_1_final.txt
```

```
ke@NGS:~/Desktop/MarkerDevelopment/data$ cat
```

```
cleaned_data/CGRL_index15_1_final.txt cleaned_data/CGRL_index15_u_final.txt | sed 's/\2$/\1/g' > raw_assembly/CGRL_index15_1_final.txt
```

```
ke@NGS:~/Desktop/MarkerDevelopment/data$ cat
```

```
cleaned_data/CGRL_index40_1_final.txt cleaned_data/CGRL_index40_u_final.txt | sed 's/\2$/\1/g' > raw_assembly/CGRL_index40_1_final.txt
```

#Copy read2 of all libraries to "raw_assembly"

```
ke@NGS:~/Desktop/MarkerDevelopment/data$ cp
```

```
cleaned_data/CGRL_index*_2_final.txt raw_assembly/
```

Commands:

#Run Trinity on 4 processors.

```
ke@NGS:~/Desktop/MarkerDevelopment/data$ 3-GenerateAssemblies trinity -a
```

```
raw_assembly/ -c 5 -e 4
```

Note: Your laptop may not be able to handle Trinity assemblies.

Output:

There are quite few intermediate files generated in

"~/Desktop/MarkerDevelopment/data/raw_assembly/CGRL_index1/".

"~/Desktop/MarkerDevelopment/data/raw_assembly/CGRL_index15/".

"~/Desktop/MarkerDevelopment/data/raw_assembly/CGRL_index40/".

#To show final trinity assemblies that are needed for annotation:

```
ke@NGS:~/Desktop/MarkerDevelopment/data$ ls raw_assembly/CGRL_index*/*.fasta
```

```

278 raw_assembly/CGRL_index15/CGRL_index15.fasta
    raw_assembly/CGRL_index1/CGRL_index1.fasta
280 raw_assembly/CGRL_index40/CGRL_index40.fasta

282 #Under "~/Desktop/MarkerDevelopment/data/" make a new folder called
    "annotation" and copy all files shown above to this folder:
284
    ke@NGS:~/Desktop/MarkerDevelopment/data$ mkdir annotation
286 ke@NGS:~/Desktop/MarkerDevelopment/data$ cp
    raw_assembly/CGRL_index*/*.fasta annotation/
288
    #check all files in folder "annotation"
290 ke@NGS:~/Desktop/MarkerDevelopment/data$ ls annotation/*
    annotation/CGRL_index15.fasta
292 annotation/CGRL_index40.fasta
    annotation/CGRL_index1.fasta
294

296 #####

298 When we did step1-3 we used a tiny fraction of the RNAseq data for the
300 purpose of quick demonstration. To better demonstrate how to use the next
    script (4-AssemblyEvaluation) let's sample some more data from each
    individual.
302
    Please do the following before you start working on step 4:
304 ke@NGS:~/Desktop/MarkerDevelopment/data$ cp
    ~/Desktop/MarkerDevelopment/associated_data/CGRL_index*.fasta
306 annotation/

308 #####

310 _____

```

312 **4-AssemblyEvaluation** (Optional): Evaluate the quality of cDNA *de novo* assemblies.
A few example of the available functions are shown here.

314

Dependencies:

316 Blat: http://hgdownload.soe.ucsc.edu/downloads.html#source_downloads

Blastall:

318 http://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE_TYPE=BlastDocs&DOC_TYPE=Download

320

Input: Trinity assemblies for all libraries stored in

322 “~/Desktop/MarkerDevelopment/data/annotation/”

324 #Display all items in “~/Desktop/MarkerDevelopment/data/annotation/”

ke@NGS:~/Desktop/MarkerDevelopment/data/annotation\$ ls

326

CGRL_index15.fasta

328

CGRL_index1.fasta

CGRL_index40.fasta

330

a. 4-AssemblyEvaluation BASIC: function “BASIC” evaluates the quality of in-target
332 assemblies by reporting basic stats: mean, median, total length, gc%, N50 etc. It also
generates a distribution of contigs by binned lengths.

334

Commands:

336 *ke@NGS:~/Desktop/MarkerDevelopment/data \$ 4-AssemblyEvaluation BASIC -a*
annotation/

338

Output:

340 # In folder “~/Desktop/MarkerDevelopment/data/annotation/”, you should get the
following output files:

342

CGRL_index15.hist

344

CGRL_index1.hist

CGRL_index40.hist

346

basic_evaluation.out

348 **Output:**

1. “XXX.hist” shows distribution of contigs by binned lengths

350

#Display first few lines of the file:

352 *ke@NGS:~/Desktop/MarkerDevelopment/data/annotation\$ head CGRL_index15.hist*

200:299 57

354

300:399 43

400:499 34

356

500:599 28

600:699 25

```

358 700:799      23
      800:899      18
360 900:999      24
      1000:1099    13
362 1100:1199    6

364 2. "basic_evaluation.out": results of assembly evaluation
      #Display first few lines of the file:
366 ke@NGS:~/Desktop/MarkerDevelopment/data/annotation$ head
      basic_evaluation.out
368
      b. 4-AssemblyEvaluation ANNOTATABLE: Calculates the percentage of the assembled
370 contigs that get a match in reference. It also calculates average percentage of
      matched bp and mismatches among the matched genes.
372
      Commands:
374 ke@NGS:~/Desktop/MarkerDevelopment/data$ 4-AssemblyEvaluation
      ANNOTATABLE -a annotation/ -b 100 -c
376 ~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.cdna.all.f
      a
378
      Output:
380 #Display results in the output file "annotatable.out":
      ke@NGS:~/Desktop/MarkerDevelopment/data/annotation$ less annotatable.out
382
      Assemblies  total matches(%)  matched bases(%)  avg similarity(%)
384 CGRL_index1    100.00 61.91 78.55
      CGRL_index15 98.00 58.38 77.36
386 CGRL_index40 96.00 68.23 78.17

388 c. 4-AssemblyEvaluation ACCURACY: The percentage of the correctly assembled
      bases estimated using the set of expressed reference transcripts
390
      Commands:
392 ke@NGS:~/Desktop/MarkerDevelopment/data$ 4-AssemblyEvaluation ACCURACY -a
      annotation/ -b 300 -c
394 ~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.pep.all.fa
396
      Output:
398 #Display results in the output file "accuracy.out":
      ke@NGS:~/Desktop/MarkerDevelopment/data/annotation$ less accuracy.out

400 Assemblies  stop codon(%)  gaps(%)
      CGRL_index1 0.000 0.000
402 CGRL_index15 0.692 0.000
      CGRL_index40 0.348 0.000

```

404 d. *4-AssemblyEvaluation CONTIGUITY*: Calculates assembly contiguity (the
 406 percentage of expressed reference transcripts covered by a single, longest
 408 assembled contig) and completeness (the percentage of expressed reference
 transcripts covered by all matched assembled contigs)

410 **Commands:**
 412 *ke@NGS:~/Desktop/MarkerDevelopment/data\$ 4-AssemblyEvaluation CONTIGUITY -*
a annotation/ -b 300 -c
 414 *~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.cdna.all.f*
a

416 ****Note:** that -b in function “CONTIGUITY” refers to the number of randomly selected
 418 sequences from the reference protein database. In functions “BASIC”,
 “ANNOTATABLE” and “ACCURACY” -b refers to the number of randomly selected
 sequences in de novo assemblies**

420 **Output:**
 422 #Display results in the output file “Contiguity.out”:
 424 *ke@NGS:~/Desktop/MarkerDevelopment/data/annotation\$ less Contiguity.out*

<i>Assemblies</i>	<i>complete(%)</i>	<i>contiguity(%)</i>
<i>CGRL_index1</i>	<i>13.46</i>	<i>11.47</i>
<i>CGRL_index15</i>	<i>23.36</i>	<i>23.36</i>
<i>GRL_index40</i>	<i>37.40</i>	<i>30.21</i>

430 _____

5-Annotation: annotate assembled contigs using a related reference protein dataset that can be found in Ensembl Genome Browser (<http://www.ensembl.org/index.html>)

Dependencies:
BLAST+:

http://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE_TYPE=BlastDocs&DOC_TYPE=Download

FrameDP: https://iant.toulouse.inra.fr/FrameDP/cgi-bin/framedp.cgi?_wb_cfg=/www/iant/FrameDP/cgi-bin/./cfg/FrameDP.cfg&_wb_session=WBUPAWho&_wb_main_menu=Download&_wb_function=Download

exonerate: <http://www.ebi.ac.uk/~guy/exonerate/index.html>

Note: this script works only if you can find a reference database from the EGB. However, if you would like to use NCBI refseq or UniProtKB/Swiss-Prot, modification of this script is needed.

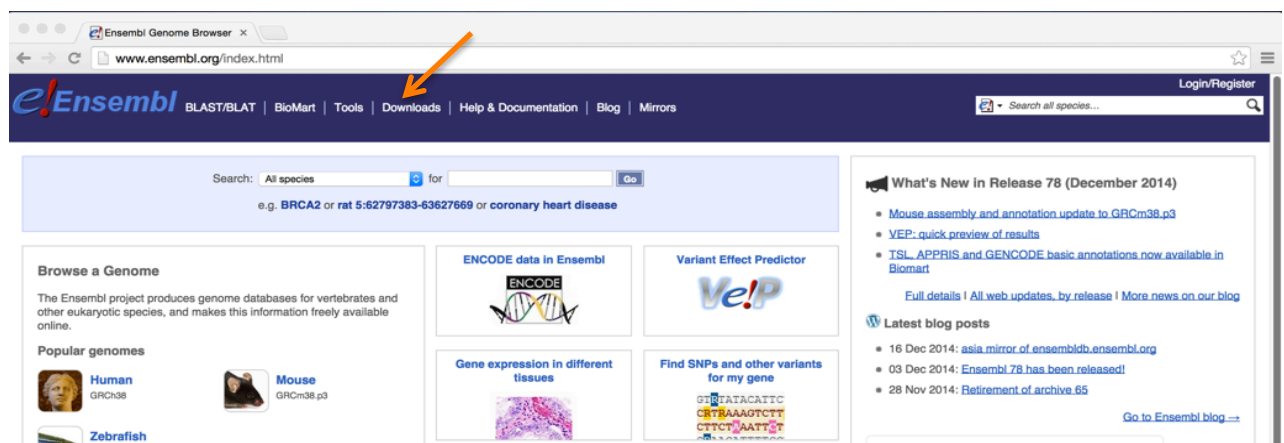
Swiss-Prot (created in 1986) is a high quality manually annotated and non-redundant protein sequence database, which brings together experimental results, computed features and scientific conclusions. UniProtKB/Swiss-Prot is now the reviewed section of the UniProt Knowledgebase.

FrameDP: Sensitive peptide detection on noisy matured sequences. A self-training integrative pipeline for predicting CDS in transcripts which can adapt itself to different levels of sequence qualities.

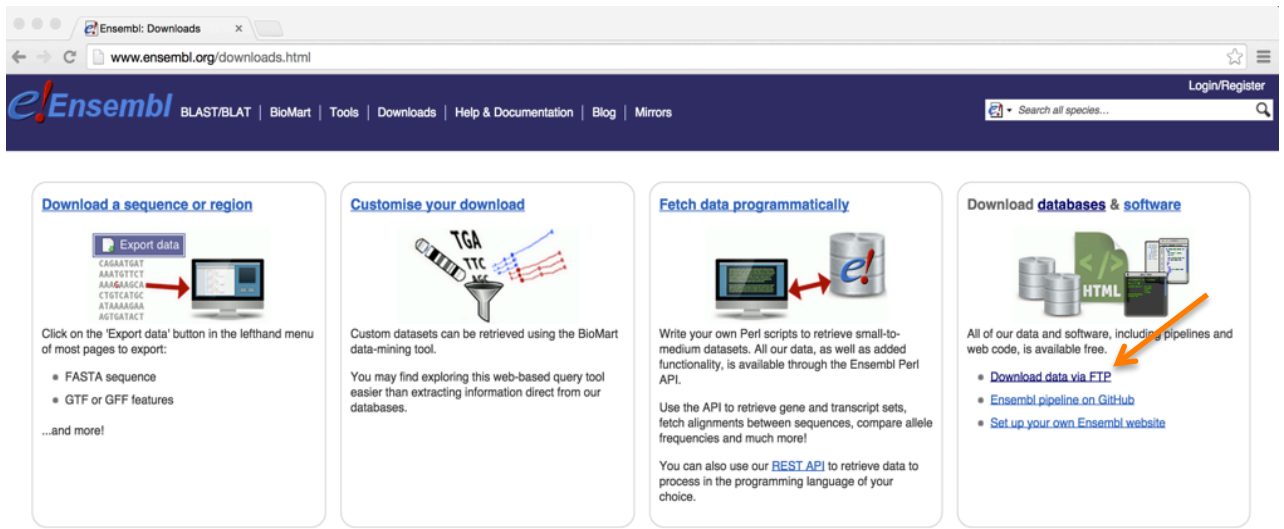
Input:

1. download a reference protein dataset from the Ensembl:

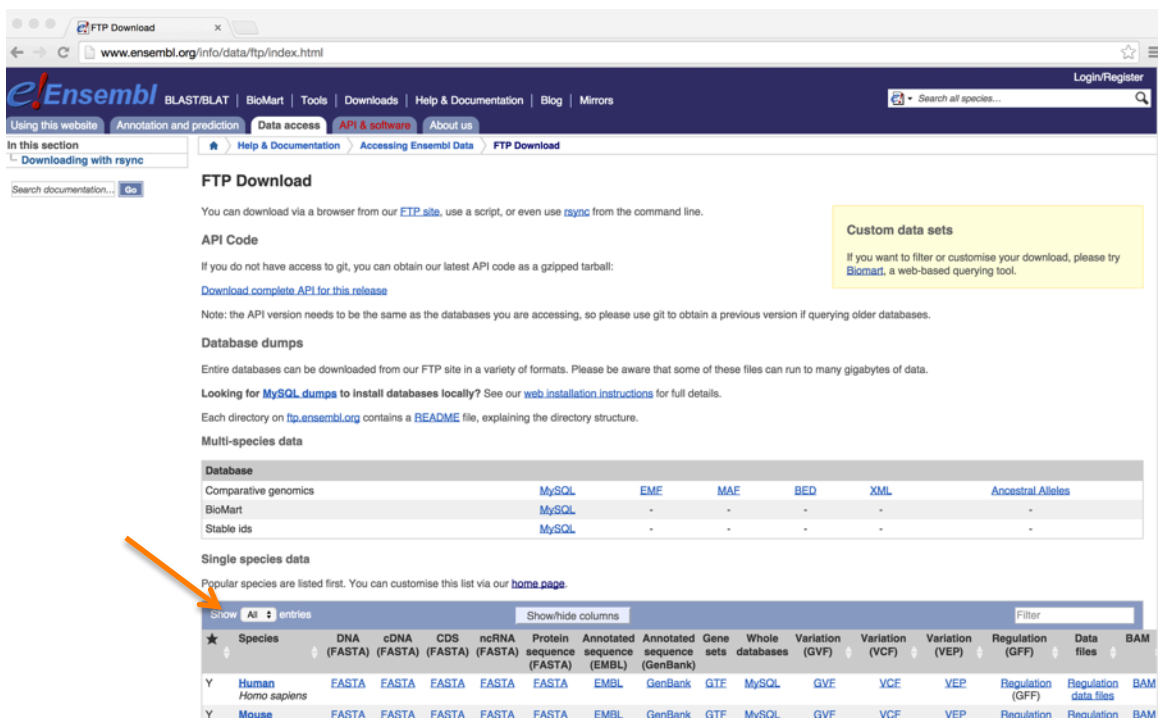
Step1. Go to the Ensembl homepage <http://www.ensembl.org/> and click on "Download" located at the top.



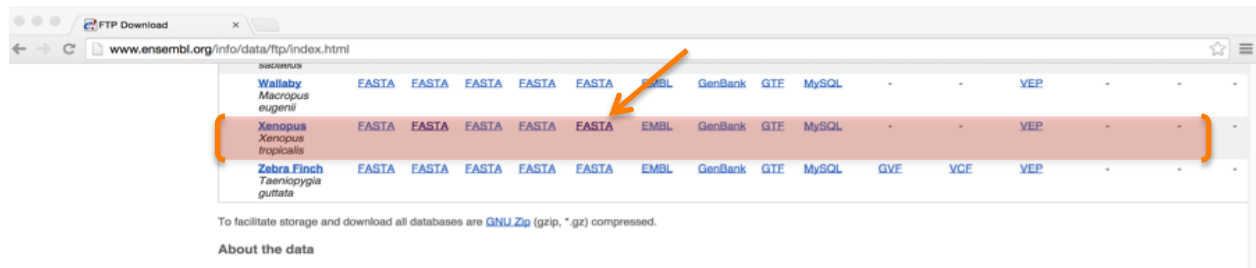
Step2. Click on “Download data via FTP” to the left of the download page.



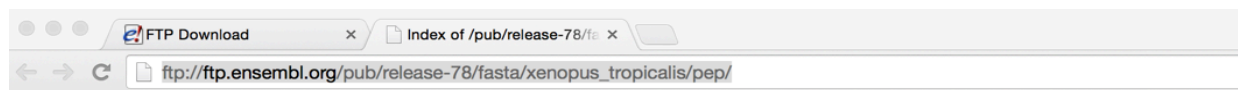
Step3. Select “All” in the “single species data” box in the FTP download page.



Step 4: Find and download the reference. Click on the FASTA link for Protein sequence. In this case we choose *Xenopus tropicalis* as the reference.



Step 5: From FTP server, download reference protein fasta “XXX.pep.all. fa.gz”



Index of /pub/release-78/fasta/xenopus_tropicalis/pep/

Name	Size	Date Modified
[parent directory]		
CHECKSUMS	130 B	11/20/14, 4:12:00 PM
README	3.0 kB	11/19/14, 1:16:00 PM
Xenopus_tropicalis.JGI_4.2.pep.abinitio.fa.gz	11.8 MB	11/19/14, 1:16:00 PM
Xenopus_tropicalis.JGI_4.2.pep.all.fa.gz	6.7 MB	11/19/14, 1:02:00 PM

Step 6: unzip the downloaded reference fasta: `gunzip Xenopus_tropicalis.JGI_4.2.pep.all. fa.gz`

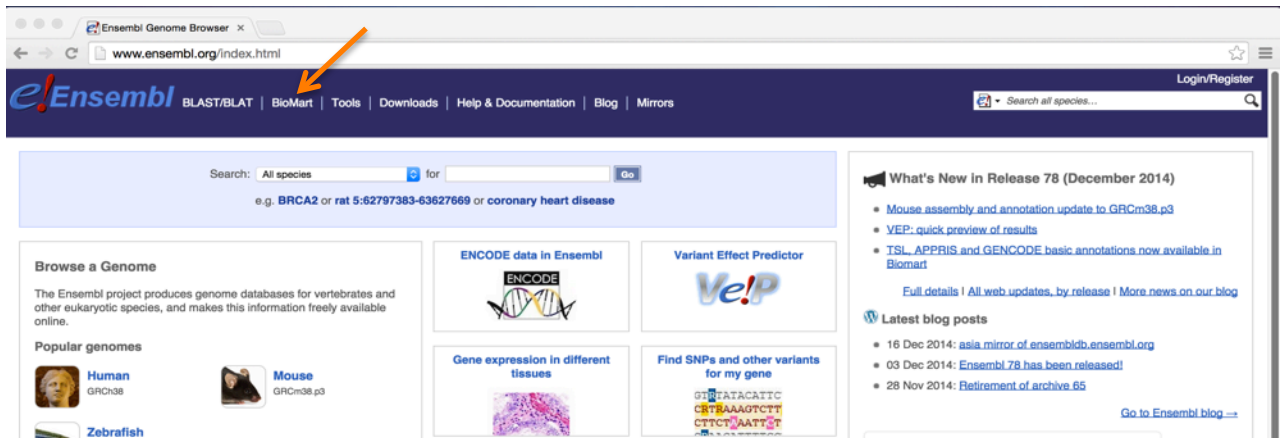
Step 7: Find and download the GTF (Gene transfer format (GTF) is a file format used to hold information about gene structure) if there is one available for the reference. In this case we can see that *Xenopus tropicalis* has a GTF so we can download it.



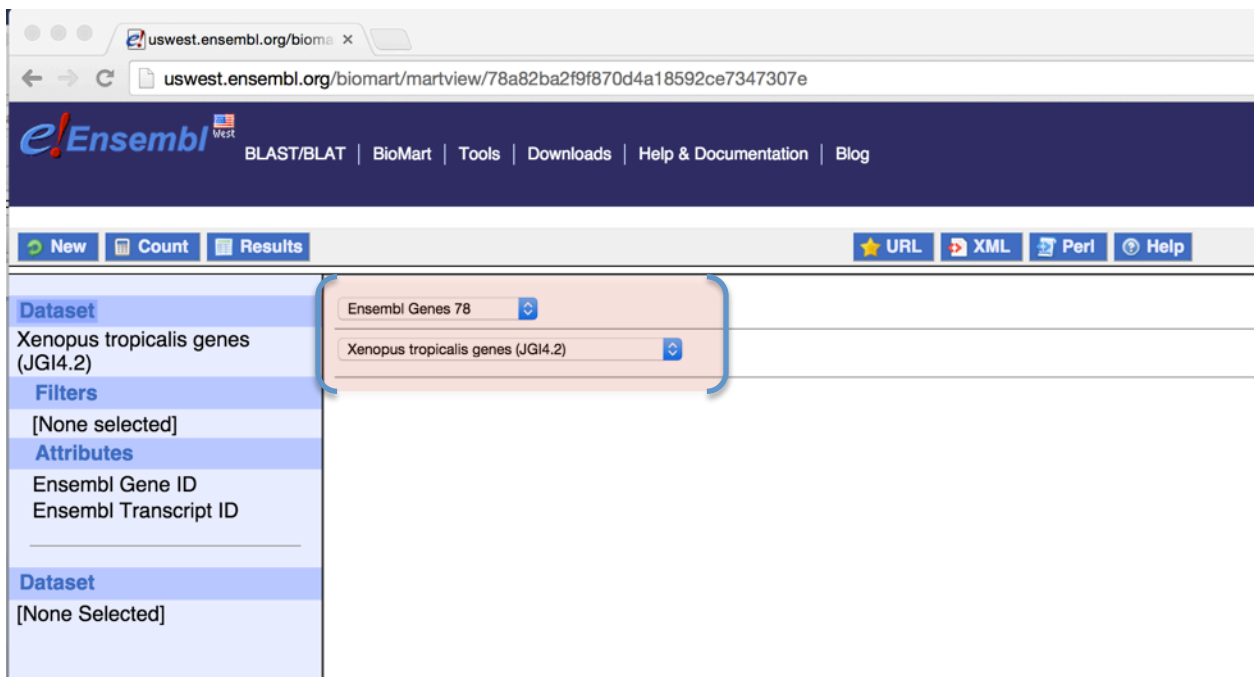
Step 8: unzip the downloaded GTF: `gunzip Xenopus_tropicalis.JGI_4.2.78.gtf.gz`

2. If GTF is not available then you can use Ensembl BioMart tool to obtain a gene annotation file for the reference. For the workshop I will show you how obtain this file from the BioMart tool even though we have downloaded a GTF for the reference.

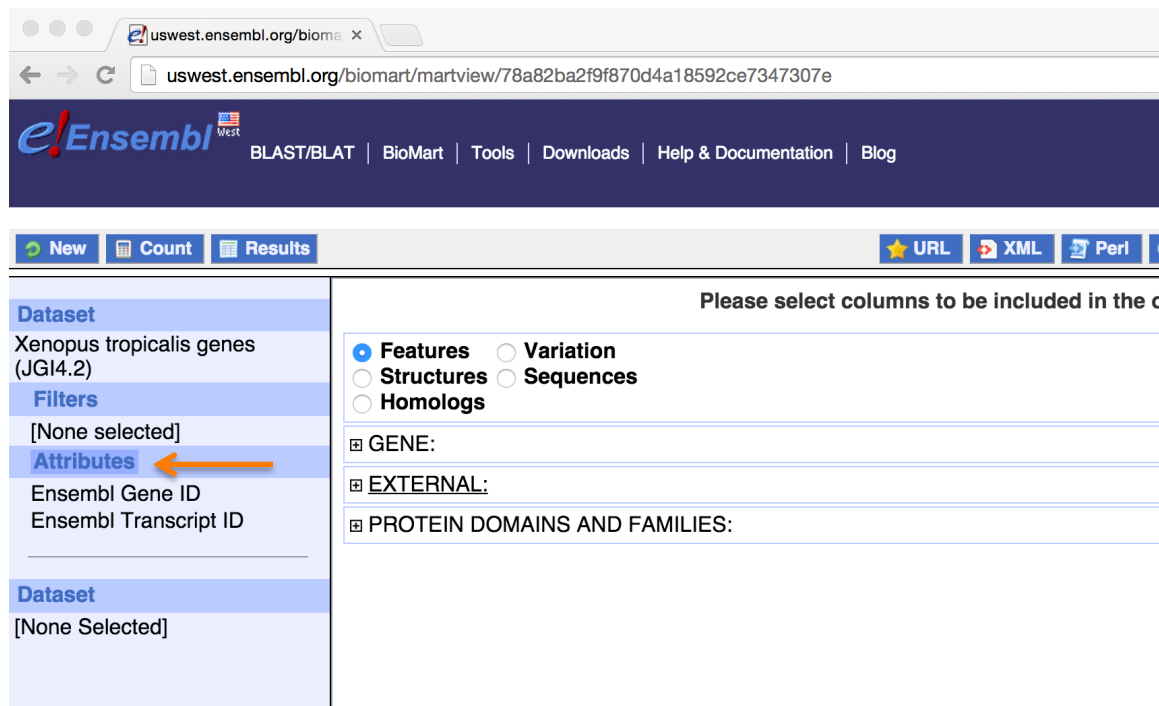
Step1. Go to the Ensembl homepage <http://www.ensembl.org/> and click on “BioMart” located at the top.



Step2. In the BioMart homepage, select “Ensembl Genes 78” and “Xenopus tropicalis genes (JGI4.2)”.

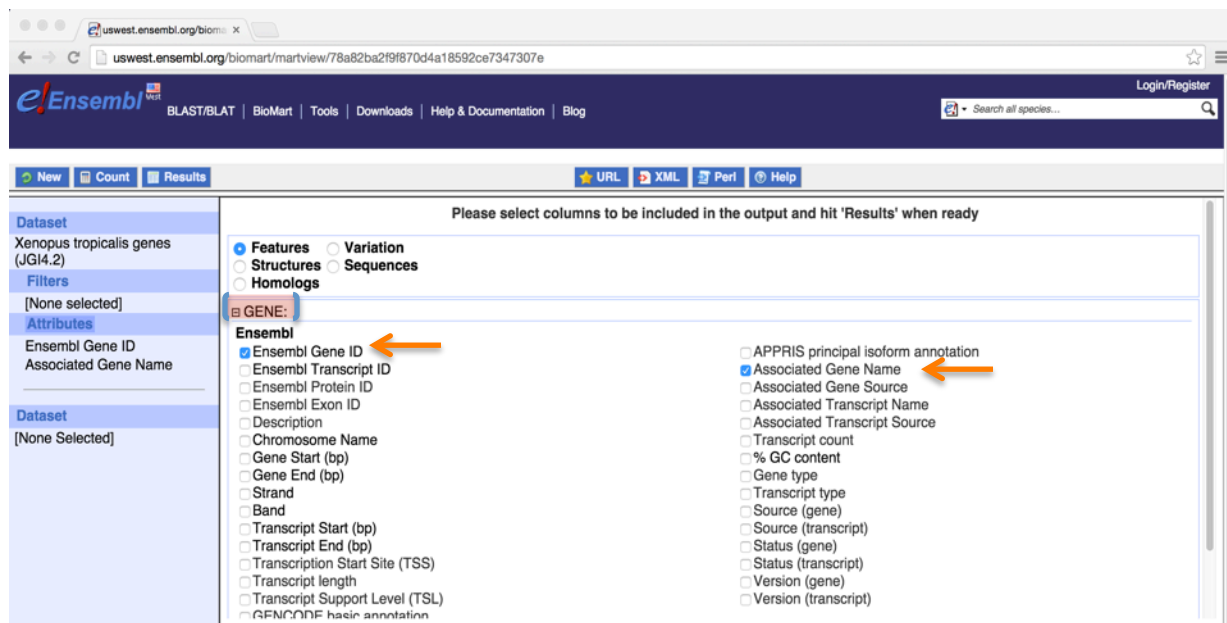


Step3. Click on “Attributes” icon to the left.



The screenshot shows the Ensembl BioMart interface. The left sidebar has a 'Dataset' section with 'Xenopus tropicalis genes (JGI4.2)' and a 'Filters' section with 'Attributes' selected (indicated by an orange arrow). The main area has a header 'Please select columns to be included in the c' and a section 'GENE:' with sub-sections 'EXTERNAL:' and 'PROTEIN DOMAINS AND FAMILIES:'. The 'Features' radio button is selected.

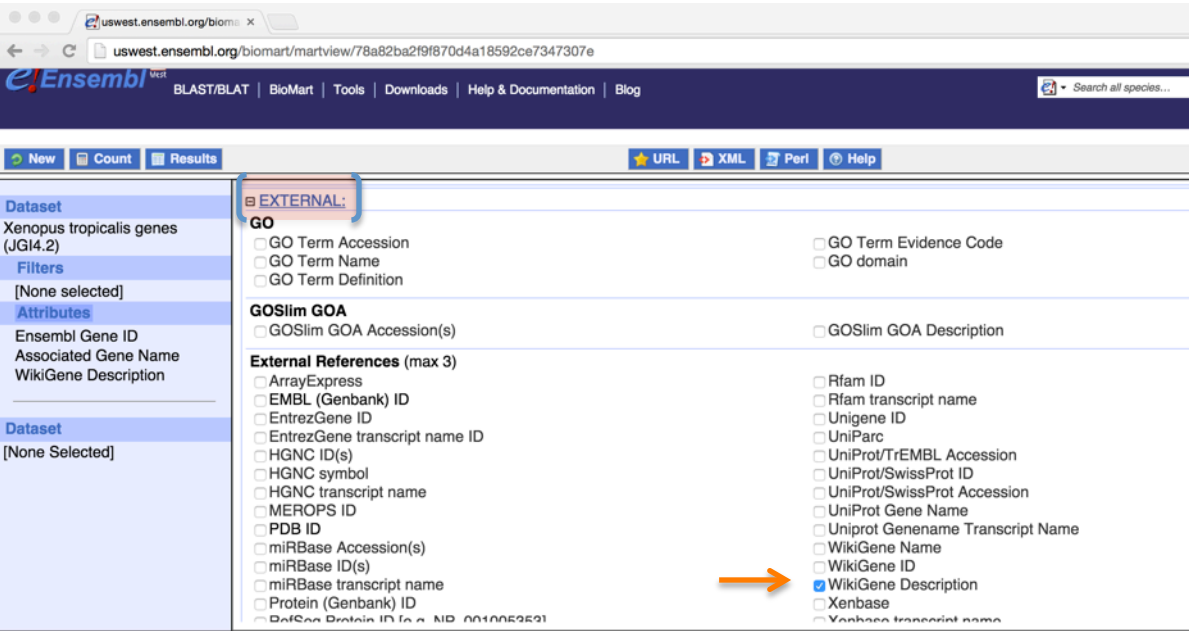
542 Step 4. Click on “GENE” to expand the manual. Check on “Ensembl Gene ID” and “Associated Gene Name”.



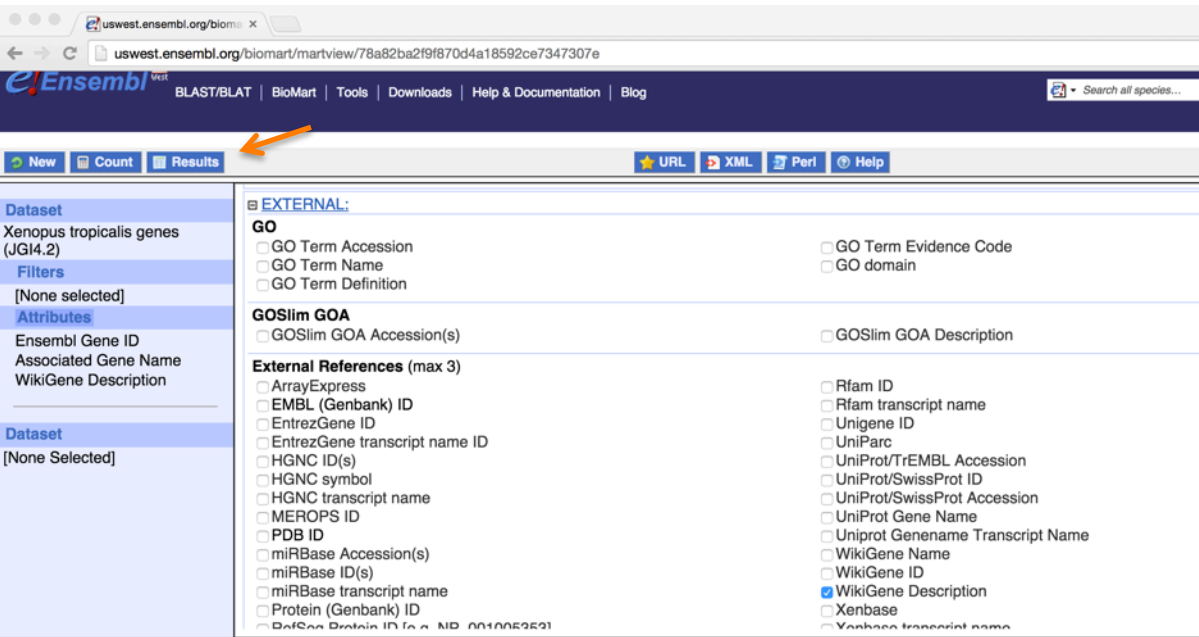
The screenshot shows the Ensembl BioMart interface. The left sidebar has a 'Dataset' section with 'Xenopus tropicalis genes (JGI4.2)' and a 'Filters' section with 'Attributes' selected. The 'GENE' section is expanded, showing a list of columns to be included in the output. The 'Ensembl Gene ID' and 'Associated Gene Name' checkboxes are checked (indicated by orange arrows). The main area has a header 'Please select columns to be included in the output and hit 'Results' when ready' and a list of columns to be included in the output.

544

546 Step 5. Scroll down the window to find “EXTERNAL”. Click on it to expand the manual. Check on “WikiGene Description”



556 Step6. Click on “Results” icon.



564

566

568

570

Step 7. To export the results, select “CSV” format and check on “Unique results only” box, and then click on “Go”.

The screenshot shows the Ensembl BioMart interface. On the left, the 'Dataset' is 'Xenopus tropicalis genes (JGI4.2)' and the 'Attributes' are 'Ensembl Gene ID', 'Associated Gene Name', and 'WikiGene Description'. The 'Export' section is set to 'all results to' 'File' in 'csv' format with 'Unique results only' checked. The 'View' section shows '10 rows as HTML' with 'Unique results only' unchecked. A table of gene data is displayed below.

Ensembl Gene ID	Associated Gene Name	WikiGene Description
ENSXETG00000002632		
ENSXETG000000026421		
ENSXETG00000008383	golt1b	golgi transport 1B
ENSXETG000000015940		
ENSXETG000000019740		
ENSXETG000000034059		coxsackievirus and adenovirus receptor homolog
ENSXETG000000016254	pafah2	
ENSXETG000000020722	selenbp1	
ENSXETG000000001197	commd7	COMM domain containing 7
ENSXETG000000034332		

Step 8. Save and rename the result to be “*Xenopus.tropicalis*_gene_name.txt”. There are three columns, separated by comma:

Ensembl Gene ID, Associated Gene Name, WikiGene Description
 ENSXETG00000008383, golt1b, golgi transport 1B
 ENSXETG000000034059, CARH, coxsackievirus and adenovirus receptor homolog
 ENSXETG000000001197, commd7, COMM domain containing 7

.....

**For this workshop, a reference protein, a GTF and the corresponding biomart gene name file are already downloaded and located in “~/Desktop/MarkerDevelopment/associated_data/”.

Input:

1. A folder that contains all trinity assemblies. These files are located in “~/Desktop/MarkerDevelopment/data/annotation/”

2. Reference protein downloaded from the ensemble:
 Xenopus_tropicalis.JGI_4.2.pep.all.fa.

3. Reference biomart gene annotation file:
 Xenopus_tropicalis_gene_name.txt

4. Reference gtf file:

```

608 Xenopus_tropicalis.JGI_4.2.78.gtf

610 Commands:
    # Run 5-Annotation without a GTF (do not execute the command during the
612 workshop, since the runs will take quite a while to finish).

614 ke@NGS:~/Desktop/MarkerDevelopment/data$ 5-Annotation -a
    ~/Desktop/MarkerDevelopment/data/annotation/ -b
616 ~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.pep.all.fa
    -d ~/Desktop/SeqCap/programs/framedp-1.2/ -f
618 ~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis_gene_name.txt -
    n xenopus -e 1

620 ##Copy the annotation results to “~/Desktop/MarkerDevelopment/data”
622 ke@NGS:~/Desktop/MarkerDevelopment/data$ scp -r
    ~/Desktop/MarkerDevelopment/associated_data/annotation/* annotation/
624

626 Output:
    For each individual trinity assembly, a new folder is generated under
    “~/Desktop/MarkerDevelopment/data/annotation/”:
628
    CGRL_index1_xenopus/
630 CGRL_index14_xenopus /
    CGRL_index40_xenopus /
632

    ##The annotated fasta files are named as “XXX_xenopus_annotated.fasta”.
634
    ke@NGS:~/Desktop/MarkerDevelopment/data/annotation$ ls
636 CGRL_index*/*annotated.fasta

638 CGRL_index15_xenopus/CGRL_index15_xenopus_annotated.fasta
    CGRL_index50_xenopus/CGRL_index50_xenopus_annotated.fasta
640 CGRL_index1_xenopus/CGRL_index1_xenopus_annotated.fasta

642 ##make a new folder “probe_design” under
    “~/Desktop/MarkerDevelopment/data”.
644 ke@NGS:~/Desktop/MarkerDevelopment/data$ mkdir probe_design

646 ##copy all the annotated fasta files to “probe_design”
    ke@NGS:~/Desktop/MarkerDevelopment/data$ cp
648 annotation/CGRL_index*/*annotated.fasta probe_design/

650 ## read and display the first few lines in the annotated fasta file:
    ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design$ head -4
652 CGRL_index15_xenopus_annotated.fasta

```

```

>contig1    gs1_ge432    ENSXETG00000014175    vwa5a NA    5e-57
654 TCTCTTACATGGACCCTTCC.....
>contig10    5u355_gs356_ge817_3u818 ENSXETG00000004176    mocs2
656 molybdenum cofactor synthesis 2 2e-82
TGTGCACAGTGTGATGTAG.....
658
For contig1: "gs1" means coding region starts at position 1. "ge432" means coding
660 region ends by position 432. No UTRs are present in this contig.
"ENSXETG00000014175" is the Ensembl gene ID obtained from Xenopus reference
662 database. "vwa5a" is the gene name. "NA" is the wiki gene description and in this
case, wiki gene description is missing. "5e-57" is e-value in the BLAST search.
664
For contig10: "5u355" means 5UTR ends by position 355. "gs356" means coding
666 region starts at position 356. "ge817" means coding region ends by position 817.
"3u818" means 3UTR starts at position 818. "ENSXETG00000004176" is the
668 Ensembl gene ID obtained from Xenopus reference database. "mocs2" is the gene
name. "molybdenum cofactor synthesis 2" is the wiki gene description. "2e-82" is e-
670 value in the BLAST search.

672 ~~~~~
Run 5-Annotation with a GTF
674
Commands:
676 ke@NGS:~/Desktop/MarkerDevelopment$ 5-Annotation -a
~/Desktop/MarkerDevelopment/data/annotation/ -b
678 ~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.pep.all.fa
-d ~/Desktop/SeqCap/programs/framedp-1.2.2/ -n xenopus -g
680 ~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.78.gtf -e
1
682
The output by using GTF is slightly different since the header doesn't have gene
684 name descriptions. For example:
>contig1    gs1_ge432    ENSXETG00000014175    vwa5a protein_coding    5e-
686 57
TCTCTTACATGGACCCTTCC.....
688
"gs1" means coding region starts at position 1. "ge432" means coding region ends by
690 position 432. No UTRs are present in this contig. "ENSXETG00000014175" is the
Ensembl gene ID obtained from Xenopus reference database. "vwa5a" is the gene
692 name. "protein_coding" is the type of the gene. "5e-57" is e-value in the BLAST
search.
694
696
698

```

700 **6-MarkerSelectionTRANS**: Find orthologous transcripts in transcriptomes from
702 different species and generate input files for probe design. It can be used when exon
identification is impossible and/or is not preferred.

704 Dependencies:

BLAST+:

706 http://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE_TYPE=BlastDocs&DOC_TYPE=Download

708 MUSCLE: <http://www.drive5.com/muscle/>

cd-hit-est: <http://weizhongli-lab.org/cd-hit/>

710

712 First of all we want to identify orthologous transcripts across transcriptomes from
different species. We will run the command “6-MarkerSelectionTRANS markers” for
this task:

714

Input:

716 All annotated transcripts located in

“~/Desktop/MarkerDevelopment/data/probe_design”

718

##

720 *ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design\$ ls*

CGRL_index15_xenopus_annotated.fasta

722 *CGRL_index40_xenopus_annotated.fasta*

CGRL_index1_xenopus_annotated.fasta

724

Make a new folder “other_files” under

726 “~/Desktop/MarkerDevelopment/data/probe_design/”.

Use one of the annotated files as a “primary” annotation file. Move the rest to a
728 folder “other_files”. In the workshop we use CGRL_index1 as the “primary”
annotation file.

730

ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design\$ mkdir other_files

732

ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design\$ mv

734 *CGRL_index15_xenopus_annotated.fasta CGRL_index40_xenopus_annotated.fasta*
other_files/

736

Commands:

738 # Run 6-MarkerSelectionTRANS markers:

740 *ke@NGS:~/Desktop/MarkerDevelopment/data\$ 6-MarkerSelectionTRANS markers -f*
probe_design/CGRL_index1_xenopus_annotated.fasta -d probe_design/other_files/ -a
1000

742

744 **Output:**

```

746 #Under "ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design/" a new
    folder called "results" was created by the script.
    ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design$ cd results/
748
    #Markers that passed all filters are stored in "marker_kept.txt". First take a How
750 many markers are kept?
    ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design/results$ wc -l
752 marker_kept.txt
    1050 marker_kept.txt
754
    ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design/results$ less -S
756 marker_kept.txt

758 Transcript_name: Ensembl Gene ID
    avgDiv: Average sequence divergence (avg. %mismatches)
760 varianceDiv: Variance of sequence divergence
    avgLength: Average length of the marker
762 avgGC: Average CG content of the marker
    div_CGRL_index15_xenopus_annotated_vs_CGRL_index1_xenopus_annotated:
764 sequence divergence between CGRL_index15 and CGRL_index1
    div_CGRL_index15_xenopus_annotated_vs_CGRL_index40_xenopus_annotated:
766 sequence divergence between CGRL_index15 and CGRL_index40
    div_CGRL_index1_xenopus_annotated_vs_CGRL_index40_xenopus_annotated:
768 sequence divergence between CGRL_index1 and CGRL_index40

770
    #Select the markers that you would like to use for probe design. In this case choose
772 the most variable 800 markers and save them in a new file "marker_final.txt"

774 ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design/results$ tail -800
    marker_kept.txt > marker_final.txt
776
    ++++++
778 Now we will use command "6-MarkerSelectionTRANS seq" to generate input fasta
    files for probe design:
780
    Input:
782 1. A final set of markers you would like to use for probe design -
    "~/Desktop/MarkerDevelopment/data/probe_design/results/marker_final.txt"
784
    2. A folder containing all trimmed transcripts in fasta format. These files were
786 created by "6-MarkerSelectionTRANS markers" and are named as XXX.final2 -
    "~/Desktop/MarkerDevelopment/data/probe_design/results/"
788

790 Commands:

```

```

# Run 6-MarkerSelectionTRANS seq:
792 ke@NGS:~/Desktop/MarkerDevelopment/data$ 6-MarkerSelectionTRANS seq -f
probe_design/results/marker_final.txt -d probe_design/results/
794 The target size for CGRL_index15_xenopus_annotated.final2 is 700532bp!
The target size for CGRL_index1_xenopus_annotated.final2 is 701541bp!
796 The target size for CGRL_index40_xenopus_annotated.final2 is 701593bp!

798 Output:
#A new folder "Probe_Design" was created by the script. cd to this folder:
800 ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design/results$ cd
Probe_Design/
802
#Three fasta sequence files contain sequences of orthologous markers are
804 generated and ready for submission for probe design:
ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design/results/Probe_Design$
806 ls *exonic_targets.txt

808 CGRL_index15_xenopus_annotated_exonic_targets.txt
CGRL_index1_xenopus_annotated_exonic_targets.txt
810 CGRL_index40_xenopus_annotated_exonic_targets.txt

812 _____
814

```

816 *6-*MarkerSelectionEXONS**: Find orthologous exons in transcriptomes from different species and generate input files for probe design.

818 Dependencies:

820 exonerate: <http://www.ebi.ac.uk/~guy/exonerate/index.html>

822 cd-hit-est: <http://weizhongli-lab.org/cd-hit/>

824 BLAST+: http://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE_TYPE=BlastDocs&DOC_TYPE=Download

826 MUSCLE: <http://www.drive5.com/muscle/>

828 We will run “6-*MarkerSelectionEXONS* exons” to identify orthologous exons in transcriptomes from each of the species. If a .gtf file is not available then we will first use a protein and genome reference to identify exons from reference species. We will then use the identified exons from the reference to identify orthologous exons from each of the transcriptomes.

832 However, if a gtf is available then I recommend first run “ParseGTF” to obtain exonic sequences from the reference and then run “6-*MarkerSelectionEXONS* exons”.

834

836 ++++++

838 First of all we assume no .gtf is available so we have to identify exons using a reference protein and reference genome.

840 **Input:**

842 1. Under “~/Desktop/MarkerDevelopment/data/” make a new folder “probe_design_exons/”:

844 *ke@NGS:~/Desktop/MarkerDevelopment/data\$ mkdir probe_design_exons/*

846 2. copy all annotated transcripts to “~/Desktop/MarkerDevelopment/data/probe_design_exons”.

848 *ke@NGS:~/Desktop/MarkerDevelopment/data\$ cp*

850 *probe_design/CGRL_index1_xenopus_annotated.fasta*

852 *probe_design/other_files/CGRL_index* probe_design_exons/*

854 *ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design_exons\$ ls*

856 *CGRL_index15_xenopus_annotated.fasta*

858 *CGRL_index40_xenopus_annotated.fasta*

860 *CGRL_index1_xenopus_annotated.fasta*

862 3. Repeat-masked reference genome
 "Xenopus_tropicalis.JGI_4.2.dna_rm.nonchromosomal.fa"

864 4. A reference protein reference "Xenopus_tropicalis.JGI_4.2.pep.all.fa";

866 Both 3 and 4 can be downloaded through Ensembl following the instruction above.

868 For this workshop these two files are located under
 "~/Desktop/MarkerDevelopment/associated_data".

870 **Command:**

872 #Run "6-MarkerSelectionEXONS exons"
 ke@NGS:~/Desktop/MarkerDevelopment/data\$ 6-MarkerSelectionEXONS exons -p
 874 '/home/ke/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.
 pep.all.fa' -g
 876 '/home/ke/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.
 dna_rm.nonchromosomal.fa' -f
 878 ~/Desktop/MarkerDevelopment/data/probe_design_exons -E 1000

880 ** "6-MarkerSelectionEXONS exons" takes very long time to run so please do not run
 it during the workshop. Let's skip this step and copy the output files directly from
 882 "associated_data":

884 ke@NGS:~/Desktop/MarkerDevelopment/data\$ cp
 ~/Desktop/MarkerDevelopment/associated_data/probe_design_exons/*.nr
 886 ~/Desktop/MarkerDevelopment/associated_data/probe_design_exons/marker_*
 probe_design_exons/
 888

890 **Output:**
 In "~/Desktop/MarkerDevelopment/data/probe_design_exons/" there are two
 892 output files that are relevant for the next step:
 1. "marker_kept.txt": Orthologous exonic markers identified in the three species
 894 2. "marker_kept_one_exon_per_gene.txt" is a subset of "marker_kept.txt"
 ,which contains randomly selected one exon per gene.

896 In both 1 and 2, annotation of each column is explained below:

898 **exon_name:** Exon ID
 900 **avgDiv:** Average sequence divergence (avg. %mismatches)
varianceDiv: Variance of sequence divergence
 902 **avgLength:** Average length of the exons
avgGC: Average CG content of the exons
 904 **div_CGRL_index15_xenopus_annotated_vs_CGRL_index1_xenopus_annotated:**
 sequence divergence between CGRL_index15 and CGRL_index1

906 **div_CGRL_index15_xenopus_annotated_vs_CGRL_index40_xenopus**
 _annotated: sequence divergence between CGRL_index15 and CGRL_index40

908 **div_CGRL_index1_xenopus_annotated_vs_CGRL_index40_xenopus_annotated:**
 sequence divergence between CGRL_index1 and CGRL_index40

910

912

914 +++++
 Now I will demonstrate how to use *6-MarkerSelectionEXONS* exons when a gtf is available.

916

918 **Command:**
 #Run "ParseGTF":

920 *ke@NGS:~/Desktop/MarkerDevelopment/data\$ ParseGTF -f*
~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.78.gtf -g
 922 *~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.dna_rm.t*
oplevel.fa -o 100 -p 1

924

926

928 **Output:**
 #Results are stored in "exons.unique" under

930 *"/home/ke/Desktop/MarkerDevelopment/associated_data/results"*

932 *#cd to "/home/ke/Desktop/MarkerDevelopment/associated_data/results"*
ke@NGS:~/Desktop/MarkerDevelopment/data\$ cd
 934 *~/Desktop/MarkerDevelopment/associated_data/results/*

936 *#display at the results*
ke@NGS:~/Desktop/MarkerDevelopment/associated_data/results\$ less -S

938 *exons.unique*

940 *#copy "exons.unique" to*
"~/Desktop/MarkerDevelopment/data/probe_design_exons/"

942 *ke@NGS:~/Desktop/MarkerDevelopment/associated_data/results\$ cp exons.unique*
~/Desktop/MarkerDevelopment/data/probe_design_exons

944 *#copy all annotated transcripts to*
 946 *"~/Desktop/MarkerDevelopment/data/probe_design_exons/"*
ke@NGS:~/Desktop/MarkerDevelopment/data\$ cp probe_design/.fasta*
 948 *probe_design_exons/*

950 **Command:**
 #run "6-MarkerSelectionEXONS exons" (do not run it in the workshop)

```

952 ke@NGS:~/Desktop/MarkerDevelopment/data$ 6-MarkerSelectionEXONS exons -p
    '/home/ke/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.
954 pep.all.fa' -g
    '/home/ke/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.
956 dna_rm.toplevel.fa' -f ~/Desktop/MarkerDevelopment/data/probe_design_exons -E
    1000

```

Output:

```

960 Same as above:

```

```

962 Now we will run command "6-MarkerSelectionEXONS seq" to generate input fasta
    files for probe design:

```

```

964

```

Input:

- 966 1. A final set of markers you would like to use for probe design. In this case we
 choose to use one exon per gene -
 968 "~/Desktop/MarkerDevelopment/data/probe_design/results/
 marker_kept_one_exon_per_gene.txt"
- 970 2. A folder containing non-redundant exonic markers in fasta format. These files
 972 were created by "6-MarkerSelectionEXONS exons" and are named as XXX_exon.fa.nr
 - "~/Desktop/MarkerDevelopment/data/probe_design_exons/"

```

974

```

Commands:

```

976 # Run 6-MarkerSelectionEXONS seq:

```

```

978 ke@NGS:~/Desktop/MarkerDevelopment/data$ 6-MarkerSelectionEXONS seq -f
    probe_design_exons/marker_kept_one_exon_per_gene.txt -d probe_design_exons/

```

```

980

```

Output:

```

982 #A new folder "Probe_Design" was created by the script. cd to this folder:

```

```

    ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design_exons$ cd

```

```

984 Probe_Design/

```

```

986 #Three fasta sequence files contain sequences of orthologous exonic markers are
    generated and ready for submission for probe design:

```

```

988

```

```

    ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design_exons/Probe_Design$ ls
990 *exonic_targets.txt

```

```

992 CGRL_index15_xenopus_annotated_exonic_targets.txt

```

```

    CGRL_index1_xenopus_annotated_exonic_targets.txt

```

```

994 CGRL_index40_xenopus_annotated_exonic_targets.txt

```