

# Pipelines for Marker Development for Transcriptome-based Exon Capture

## *Part II Population Genomics*

February 17, 2015

Contributors: Sonal Singhal and Ke Bi

For questions or to report bugs, please contact Ke Bi ([kebi@berkeley.edu](mailto: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, Linderth 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, Linderth 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/MarkerDevelopmentPopGen>

---

Scripts included in this pipeline:

[1-PreCleanup](#)

[2-ScrubReads](#)

[3-GenerateAssemblies](#)

[4-AssemblyEvaluation](#)

[5-Annotation](#)

6-ProcessAnnotation

7-MiningTranscripts

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

\*\*In transcriptome-based exon captures, marker development for population

genomic projects needs RNAseq from multiple tissue types from one individual  
sample. For this workshop, we selected cDNA from three tissues from one individual  
frog and the libraries are named as "CGRL\_index1", "CGRL\_index15", and  
"CGRL\_index40", respectively.

---

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

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

62 **Input:**  
 Raw sequence data files are grouped and saved in folders named by their sample  
 64 IDs. For instance, three libraries (CGRL\_index1, CGRL\_index15, CGRL\_index40) are  
 saved under “/home/ke/Desktop/SeqCap/data/rawdata/library/”. Compressed  
 66 fastq sequence files are saved in each of these folders.

68 Fastq files use the following naming scheme:  
 <sample name>\_<barcode sequence>\_L<lane (0-padded to 3 digits)>\_R<read  
 70 number>\_<set number (0-padded to 3 digits)>.fastq.gz

72 For example, in “CGRL\_index15\_CGACCTG\_L006\_R1\_001.fastq.gz”:  
 sample name: CGRL\_index15  
 74 barcode sequence: CGACCTG  
 lane (0-padded to 3 digits): 006  
 76 read number: 1  
 set number (0-padded to 3 digits): 001

78 #Make a new folder called “raw” under  
 80 “~/Desktop/MarkerDevelopment/data/rawdata”:  
*ke@NGS:~/Desktop/MarkerDevelopment/data/rawdata\$ mkdir raw*

82 #Copy all these compressed fastq files from each folder (CGRL\_index1,  
 84 CGRL\_index15, CGRL\_index40) to “raw”:  
*ke@NGS:~/Desktop/MarkerDevelopment/data/rawdata\$ cp*  
 86 *library/CGRL\_index\*/\*.gz raw/*

88 #Check data files in “raw”:  
*ke@NGS:~/Desktop/MarkerDevelopment/data/rawdata\$ ls raw/\**  
 90 *CGRL\_index15\_CGACCTG\_L006\_R1\_001.fastq.gz*  
*CGRL\_index15\_CGACCTG\_L006\_R2\_001.fastq.gz*  
 92 *CGRL\_index1\_TCGCAGG\_L006\_R1\_001.fastq.gz*  
*CGRL\_index1\_TCGCAGG\_L006\_R2\_001.fastq.gz*  
 94 *CGRL\_index40\_TTCGCAA\_L006\_R1\_001.fastq.gz*  
*CGRL\_index40\_TTCGCAA\_L006\_R2\_001.fastq.gz*  
 96

```

98  Commands:
    #cd to the working directory:
100  ke@NGS:~/Desktop/MarkerDevelopment/data/rawdata$ cd ..

102  #run 1-PreCleanup with fastq evaluation
    ke@NGS:~/Desktop/MarkerDevelopment/data$ 1-PreCleanup
104  ~/Desktop/MarkerDevelopment/data/rawdata/raw/ fastqc

106  ~/Desktop/MarkerDevelopment/data
Output:
108  Three new folders will be created under
    "~/Desktop/MarkerDevelopment/data/rawdata/raw/":
110  "pre-clean"
    "combined"
112  "pre-clean/evaluation"

114  - Folder "pre-clean" contains reformatted raw fastq reads.
    CGRL_index1_R1.fq
116  CGRL_index1_R2.fq
    CGRL_index15_R1.fq
118  CGRL_index15_R2.fq
    CGRL_index40_R1.fq
120  CGRL_index40_R2.fq

122  - Folder "combined" contains merged, compressed, fastq data files (not used by the
    following pipeline).
124  CGRL_index1_TCGCAGG_L006_R1.fastq.gz
    CGRL_index1_TCGCAGG_L006_R2.fastq.gz
126  CGRL_index15_CGACCTG_L006_R1.fastq.gz
    CGRL_index15_CGACCTG_L006_R2.fastq.gz
128  CGRL_index40_TTCGCAA_L006_R1.fastq.gz
    CGRL_index40_TTCGCAA_L006_R2.fastq.gz
130
    - Folder "evaluation" contains fastQC results for each data file.
132  CGRL_index1_R1.fq_fastqc/
    CGRL_index1_R2.fq_fastqc/
134  CGRL_index15_R1.fq_fastqc/
    CGRL_index15_R2.fq_fastqc/
136  CGRL_index40_R1.fq_fastqc/
    CGRL_index40_R2.fq_fastqc/
138

```

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

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

152 **Input:**

154 1. Reformatted fastq files created by [1-PreCleanup](#):  
 #Check the raw data files:  
 156 *ke@NGS:~/Desktop/MarkerDevelopment/data/rawdata/raw/pre-clean\$ ls \*.fq*  
*CGRL\_index1\_R1.fq*  
 158 *CGRL\_index1\_R2.fq*  
*CGRL\_index15\_R1.fq*  
 160 *CGRL\_index15\_R2.fq*  
*CGRL\_index40\_R1.fq*  
 162 *CGRL\_index40\_R2.fq*

164 2. A fasta file that contains adapter sequences:  
 #Check the format of adapter sequence file:  
 166 *ke@NGS:~/Desktop/SeqCap/denovoTargetCapture/associated\_files \$ less -S*  
*Adapters.fasta*  
 168 *>P7\_index1*  
*CAAGCAGAAGACGGCATACGAGATcctgcgaGTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT*  
 170 *>P7\_index2*  
*CAAGCAGAAGACGGCATACGAGATtgcagagGTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT*  
 172 *.....*  
*>P5\_index1*  
 174 *AATGATACGGCGACCACCGAGATCTACACcctgcgaACACTCTTCCCTACACGACGCTCTTCCGATCT*  
*>P5\_index2*  
 176 *AATGATACGGCGACCACCGAGATCTACACtgcagagACACTCTTCCCTACACGACGCTCTTCCGATCT*  
 178 *.....*

Note: The header of each adapter sequence has to be named strictly as “**P7\_indexN**”  
 180 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.

182 3. Library info file (Tab-delimited txt file):  
 184 #Check the format of Library info file:  
*ke@NGS:~/Desktop/SeqCap/denovoTargetCapture/associated\_files \$ less -S libInfo.txt*  
 186

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

```

188  CGRL_index1  1
      CGRL_index15  15
190  CGRL_index40  40

```

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

194 4. Contaminant file:  
*Escherichia coli* (+ human + other genome resources if desired) genome in fasta  
196 format.

This file (e\_coli\_K12.fasta) is saved in  
198 “~/Desktop/SeqCap/denovoTargetCapture/associated\_files/ecoli/”

200

#### **Commands:**

```

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

```

```

206 #Run 2-ScrubReads:
      ke@NGS:~/Desktop/MarkerDevelopment/data$ 2-ScrubReads -f
208 ~/Desktop/MarkerDevelopment/data/rawdata/raw/pre-clean/ -o
      ~/Desktop/MarkerDevelopment/data/cleaned_data/ -a
210 ~/Desktop/SeqCap/denovoTargetCapture/associated_files/Adapters.fasta -b
      ~/Desktop/SeqCap/denovoTargetCapture/associated_files/libInfo.txt -t
212 /home/ke/Desktop/SeqCap/programs/Trimmomatic-0.32/trimmomatic-0.32.jar -c
      ~/Desktop/SeqCap/denovoTargetCapture/associated_files/ecoli/e_coli_K12.fasta -e
214 200 -m 15 -z

```

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

218

#### **Output:**

220 1. In “~/Desktop/MarkerDevelopment/data/cleaned\_data/”, six .txt files per  
library are produced:

```

222   For example for library CGRL_index1, the six files are:
      CGRL_index1_1_final.txt (left reads)
224   CGRL_index1_2_final.txt (right reads)
      CGRL_index1_u_final.txt (merged or unpaired reads)
226   CGRL_index1.contam.out (headers of reads aligned to bacteria)
      CGRL_index1.duplicates.out (headers of duplicated reads)
228   CGRL_index1.lowComplexity.out (headers of low complexity reads)

```

230 2. In “~/Desktop/MarkerDevelopment/data/cleaned\_data/evaluation/”, you can  
find fastQC results for cleaned reads from each library.

232

```

234  *3-GenerateAssemblies*: Assemble multi-tissue RNAseq data using Trinity.
236  Dependencies:
237  Trinity http://trinityrnaseq.sourceforge.net
238
239  Input:
240  For each library, concatenate cleaned forward reads (XXX_1_final.txt) and unpaired
241  reads (XXX_u_final.txt) and name the resulting read data file as XXX_1_final.txt.
242
243  #Make a new folder called "raw_assembly" under
244  "~/Desktop/MarkerDevelopment/data/":
245  ke@NGS:~/Desktop/MarkerDevelopment/data$ mkdir raw_assembly_pop
246
247  #Concatenate cleaned forward reads and unpaired reads and save them in
248  "raw_assembly_pop/":
249  ke@NGS:~/Desktop/MarkerDevelopment/data$ cat cleaned_data/*_1_final.txt
250  cleaned_data/*_u_final.txt | sed 's/\2$/\1/g' >
251  raw_assembly_pop/combined_1_final.txt
252
253  ke@NGS:~/Desktop/MarkerDevelopment/data$ cat cleaned_data/*_2_final.txt >
254  raw_assembly_pop/combined_2_final.txt
255
256  #Concatenated files "combined_1_final.txt" and "combined_2_final.txt" are the input
257  files for trinity assembly.
258
259
260
261  Commands:
262  #Run Trinity on 4 processors.
263  ke@NGS:~/Desktop/MarkerDevelopment/data$ 3-GenerateAssemblies trinity -a
264  raw_assembly_pop/ -c 5 -e 4
265
266  Note: Your laptop may not be able to handle memory intensive Trinity assemblies.
267
268  Output:
269  #The resulting trinity assembly is named "combined.fasta" in
270  "~/Desktop/MarkerDevelopment/data/raw_assembly_pop/combined/".
271
272  #Under "~/Desktop/MarkerDevelopment/data/", make a new folder called
273  "annotation_pop" and "combined.fasta" shown above to this folder:
274
275  ke@NGS:~/Desktop/MarkerDevelopment/data$ mkdir annotation_pop
276  ke@NGS:~/Desktop/MarkerDevelopment/data$ cp
277  raw_assembly_pop/combined/combined.fasta annotation_pop/
278

```

---

280

#####

282

284 **When we did step1-3 we used a tiny fraction of the RNAseq data for the**  
286 **purpose of quick demonstration. To better demonstrate how to use the next**  
**script (4-AssemblyEvaluation) let's sample some more data from each**  
**individual.**

288

**Please do the following before you start working on step 4:**

290

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

292

**~/Desktop/MarkerDevelopment/associated\_data/combined.fasta**  
**annotation/**

294

#####

296

---

298

298



300 *\*4-AssemblyEvaluation\** (Optional): Evaluate the quality of RNAseq data *de novo* assemblies. A few example functions are shown here.

302 Dependencies:  
 303 Blat: [http://hgdownload.soe.ucsc.edu/downloads.html#source\\_downloads](http://hgdownload.soe.ucsc.edu/downloads.html#source_downloads)  
 304 Blastall:  
 305 [http://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE\\_TYPE=BlastDocs&DOC\\_TYPE=Download](http://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE_TYPE=BlastDocs&DOC_TYPE=Download)  
 306 ad

308 **Input:** A trinity assembly “combined.fasta” stored in  
 309 “~/Desktop/MarkerDevelopment/data/annotation\_pop/”

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

314 **Commands:**

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

318 **Output:**

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

321 *combined.hist*  
 322 *basic\_evaluation.out*

326 **Output:**

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

328 #Display first few lines of the file:

329 *ke@NGS:~/Desktop/MarkerDevelopment/data/annotation\_pop\$ head combined.hist*

330	100:199	2
331	200:299	865
332	300:399	619
333	400:499	483
334	500:599	426
335	600:699	376
336	700:799	350
337	800:899	318
338	900:999	284
339	1000:1099	269

340

341 2. “basic\_evaluation.out”: results of assembly evaluation

342 #Display first few lines of the file:

343 *ke@NGS:~/Desktop/MarkerDevelopment/data/annotation\_pop\$ head*

*basic\_evaluation.out*

b. *4-AssemblyEvaluation ANNOTATABLE*: Calculates the percentage of the assembled contigs that get a match in reference. It also calculates average percentage of matched bp and mismatches among the matched genes.

**Commands:**

```
ke@NGS:~/Desktop/MarkerDevelopment/data$ 4-AssemblyEvaluation  
ANNOTATABLE -a annotation_pop/ -b 100 -c  
~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.cdna.all.f  
a
```

**Output:**

```
#Display results in the output file "annotatable.out":  
ke@NGS:~/Desktop/MarkerDevelopment/data/annotation_pop$ less annotatable.out
```

<i>Assemblies</i>	<i>total matches(%)</i>	<i>matched bases(%)</i>	<i>avg similarity(%)</i>
<i>combined</i>	<i>96.00</i>	<i>62.53</i>	<i>77.70</i>

c. *4-AssemblyEvaluation ACCURACY*: The percentage of the correctly assembled bases estimated using the set of expressed reference transcripts

**Commands:**

```
ke@NGS:~/Desktop/MarkerDevelopment/data$ 4-AssemblyEvaluation ACCURACY -a  
annotation_pop/ -b 300 -c  
~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.pep.all.fa
```

**Output:**

```
#Display results in the output file "accuracy.out":  
ke@NGS:~/Desktop/MarkerDevelopment/data/annotation_pop$ less accuracy.out
```

<i>Assemblies</i>	<i>stop codon(%)</i>	<i>gaps(%)</i>
<i>combined</i>	<i>0.355</i>	<i>0.000</i>

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

**Commands:**

```
ke@NGS:~/Desktop/MarkerDevelopment/data$ 4-AssemblyEvaluation CONTIGUITY -  
a annotation_pop/ -b 300 -c  
~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.cdna.all.f  
a
```

**\*\*Note:** that -b in function "CONTIGUITY" refers to the number of randomly selected

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\*\*

**Output:**

#Display results in the output file "Contiguity.out":  
ke@NGS:~/Desktop/MarkerDevelopment/data/annotation\_pop\$ less Contiguity.out

Assemblies	complete(%)	contiguity(%)
combined	39.86	31.63

**\*5-Annotation\*:** annotate assembled transcripts using one or multiple reference protein datasets that can be found in Ensembl Genome Browser (<http://www.ensembl.org/index.html>).

Using multiple references may be helpful when there is not a closely related reference genome available.

For this workshop we use two references, Anole lizard (*Anolis carolinensis*) and clawed frog (*Xenopus tropicalis*), to annotate the assemblies (from a frog species). The example below shows how to download protein reference and corresponding annotation files of *Xenopus tropicalis* From the Ensembl Genome Browser.

Dependencies:

BLAST+:

[http://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE\\_TYPE=BlastDocs&DOC\\_TYPE=Download](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](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.

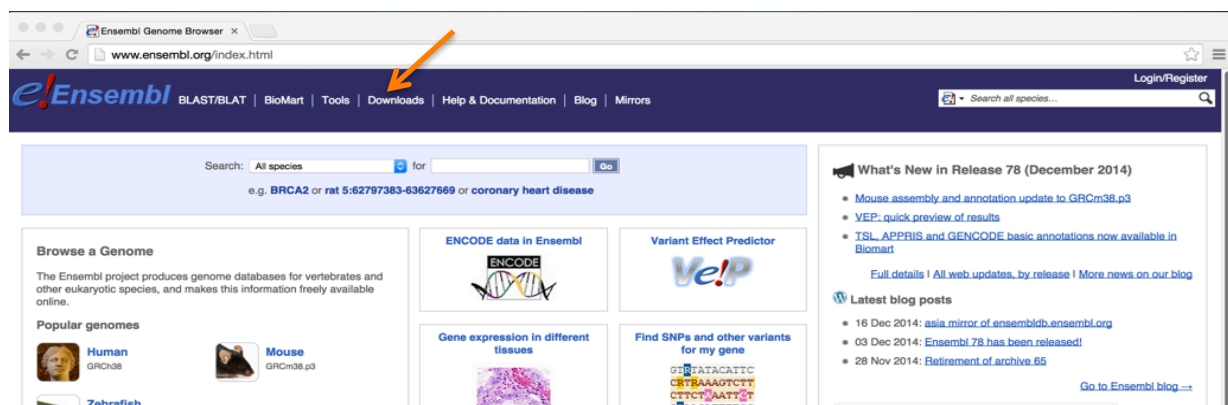
**\*\* Swiss-Prot** (created in 1986) is a high quality manually annotated and non-redundant protein sequence database, which brings together experimental results, computed features and scientific conclusions. UniProtKB/Swiss-Prot is now the reviewed section of the UniProt Knowledgebase.

**\*\* 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.

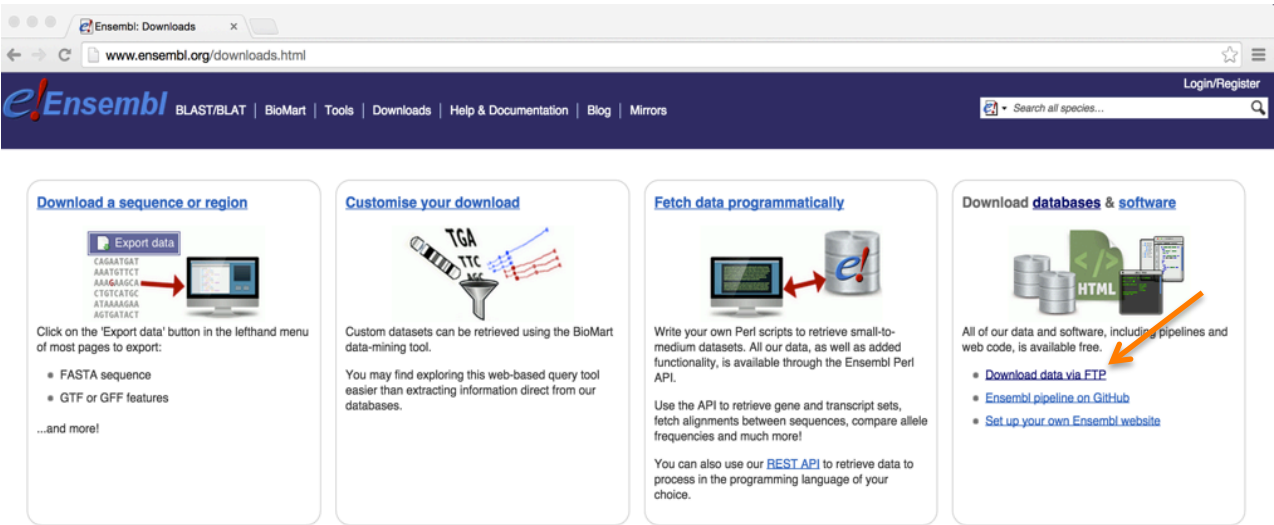
### Input:

1. download a reference protein dataset from the Ensembl:

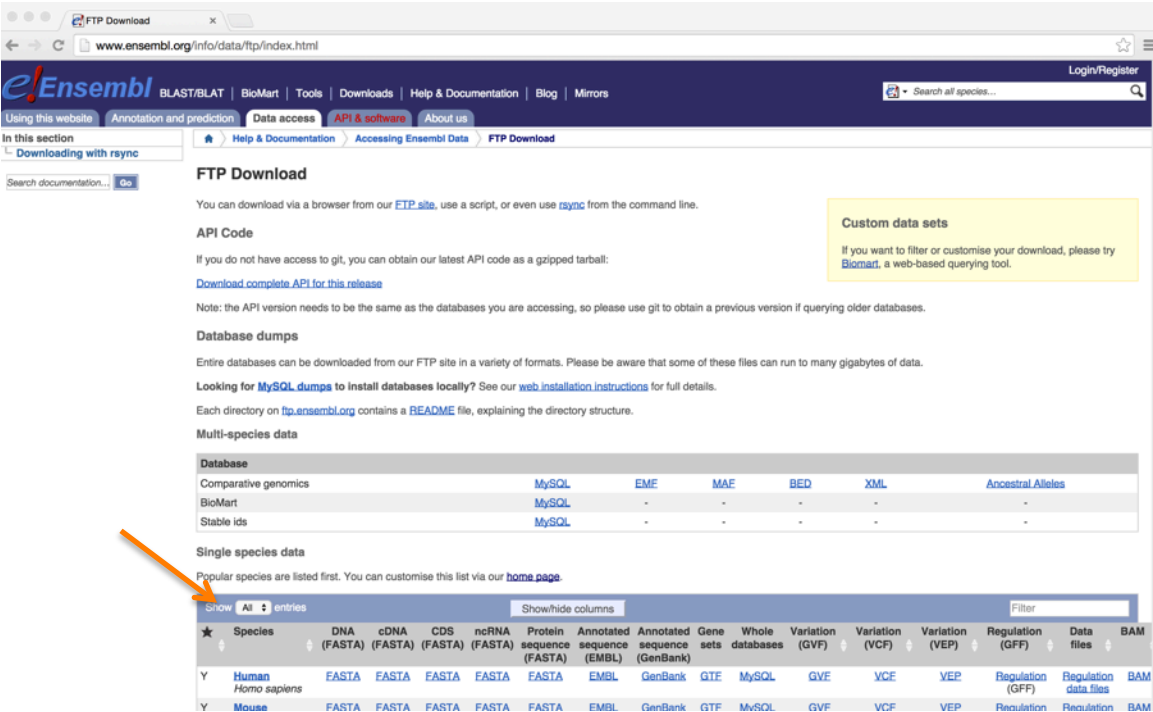
Step1. Go to the Ensembl homepage <http://www.ensembl.org/> and click on "Download" located at the top.



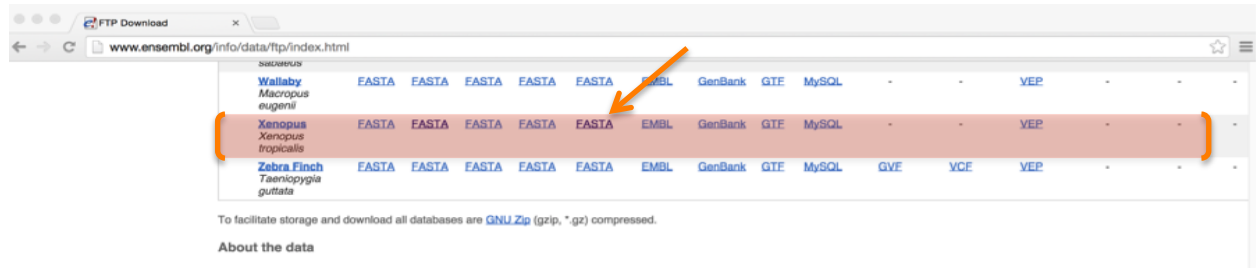
Step2. Click on “Download data via FTP” to the left of the download page.



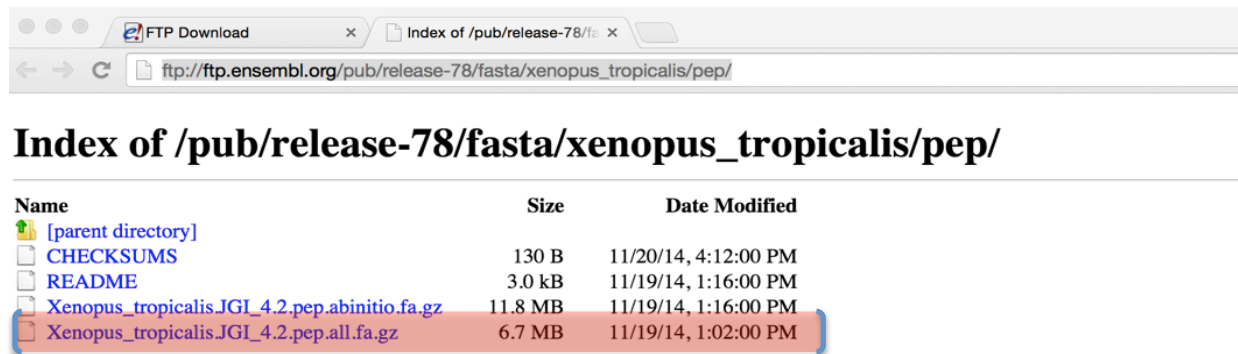
Step3. Select “All” in the “single species data” box in the FTP download page.



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”



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

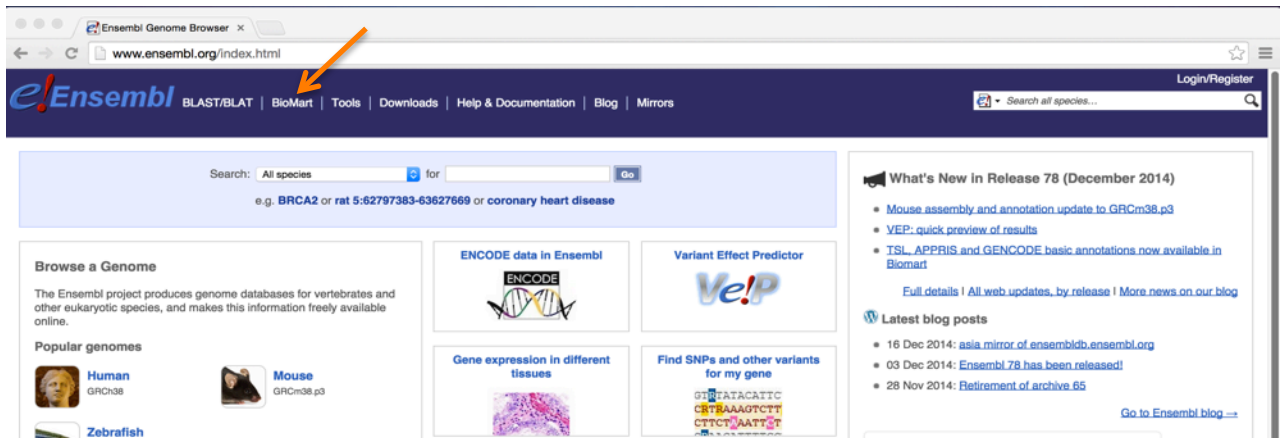
Step 7: Find and download the GTF (Gene transfer format (GTF) is a file format used to hold information about gene structure) if there is one available for the reference. In this case we can see that *Xenopus tropicalis* has a GTF so we can download it.



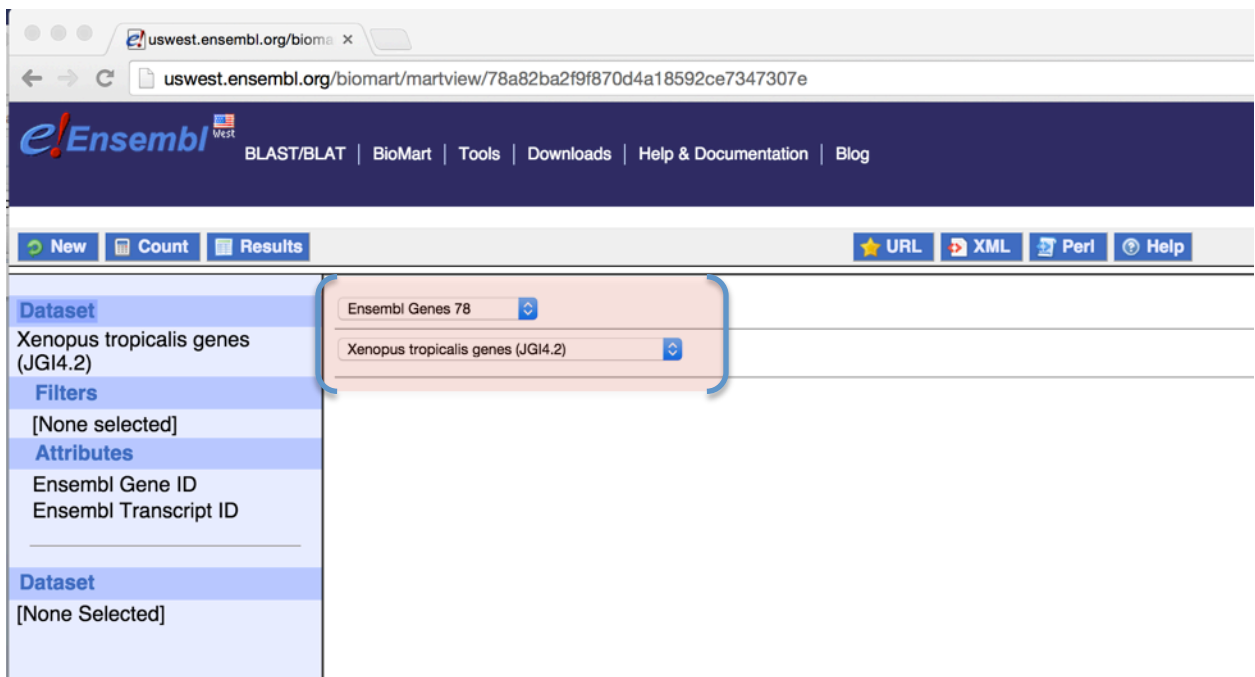
Step 8: unzip the downloaded GTF: `gunzip Anolis_carolinensis.AnoCar2.0.78.gtf.gz`

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.

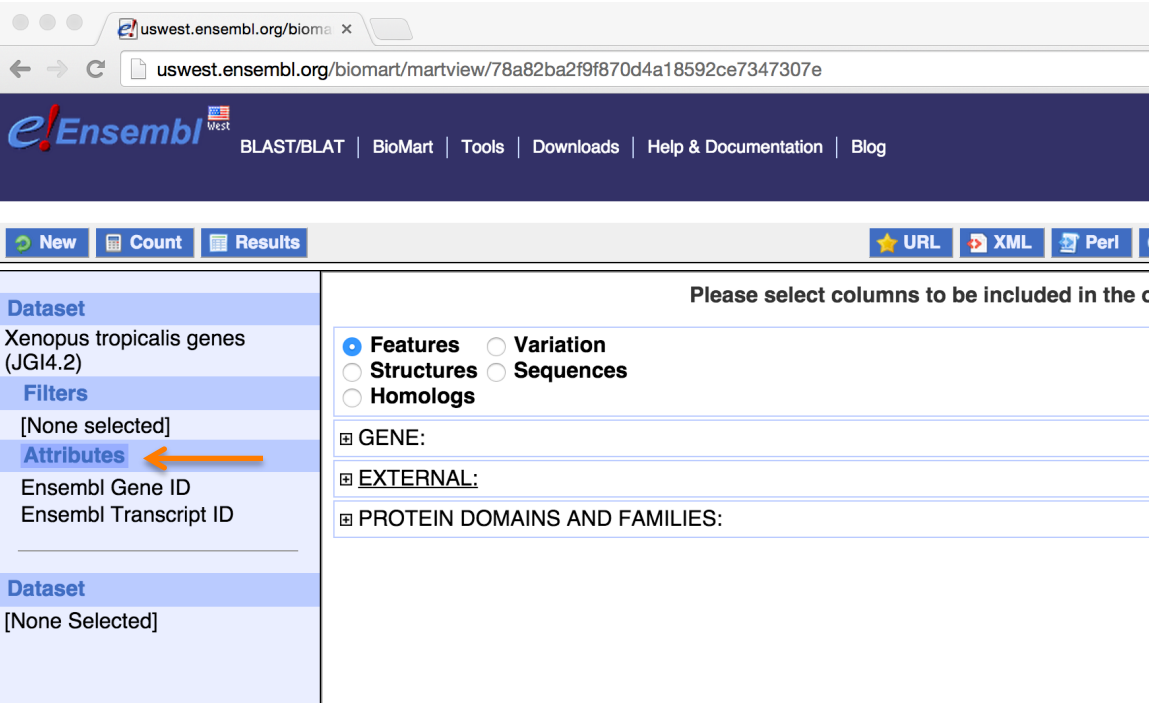
Step1. Go to the Ensembl homepage <http://www.ensembl.org/> and click on “BioMart” located at the top.



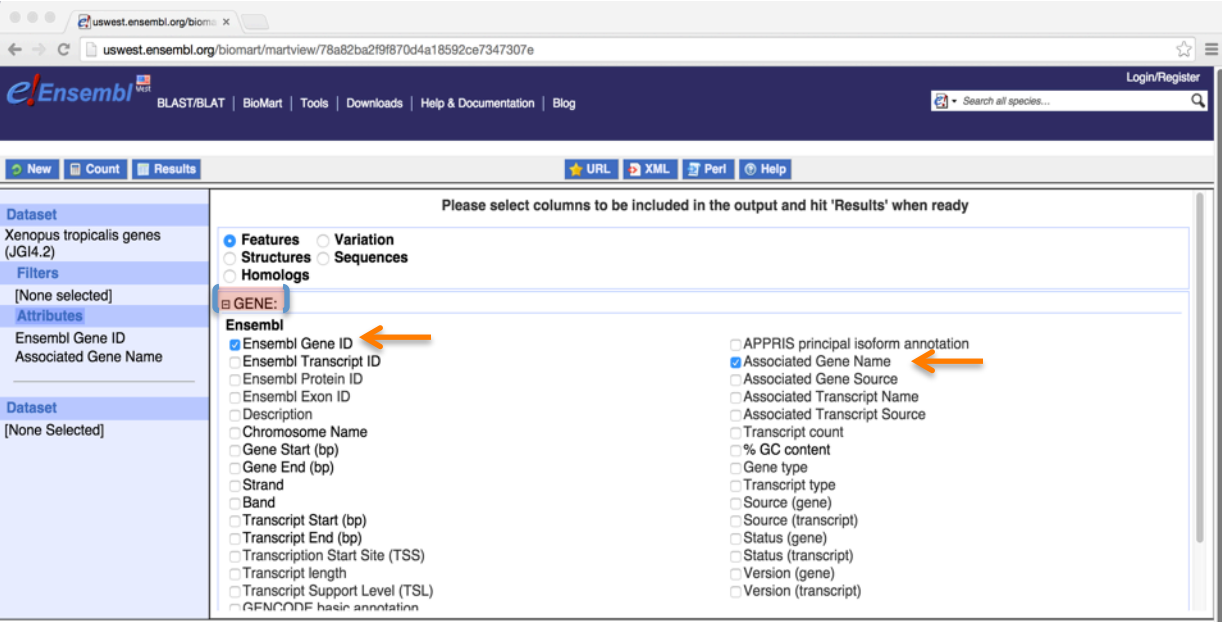
Step2. In the BioMart homepage, select “Ensembl Genes 78” and “Xenopus tropicalis genes (JGI4.2)”.



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



516  
518 Step 4. Click on “GENE” to expand the manual. Check on “Ensembl Gene ID” and “Associated Gene Name”.

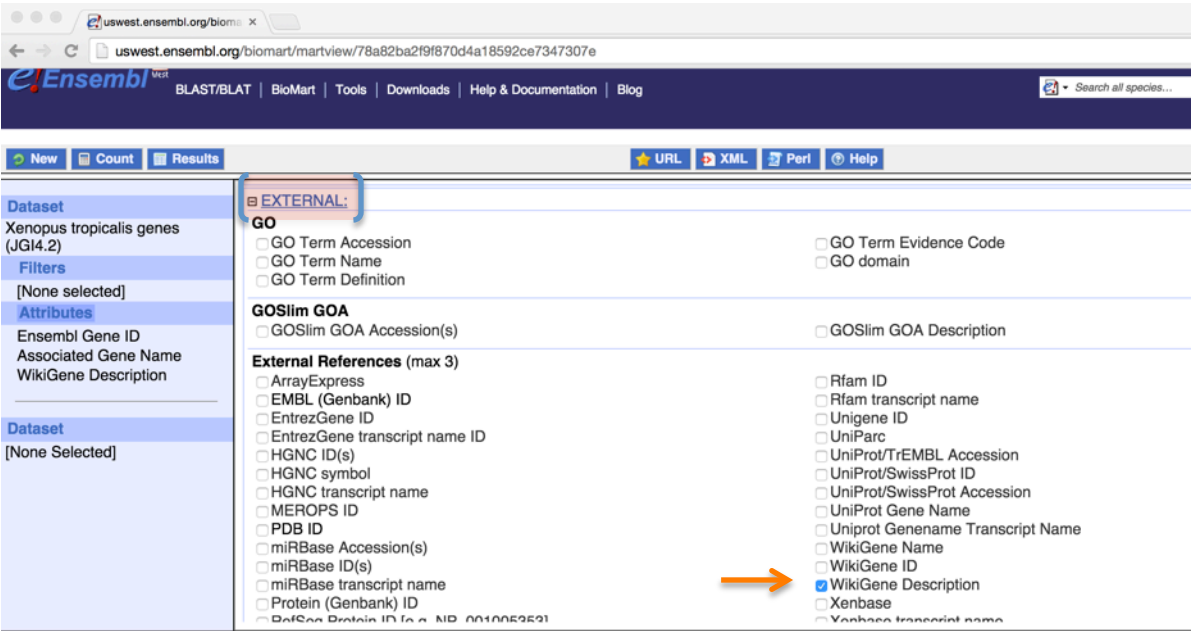


520

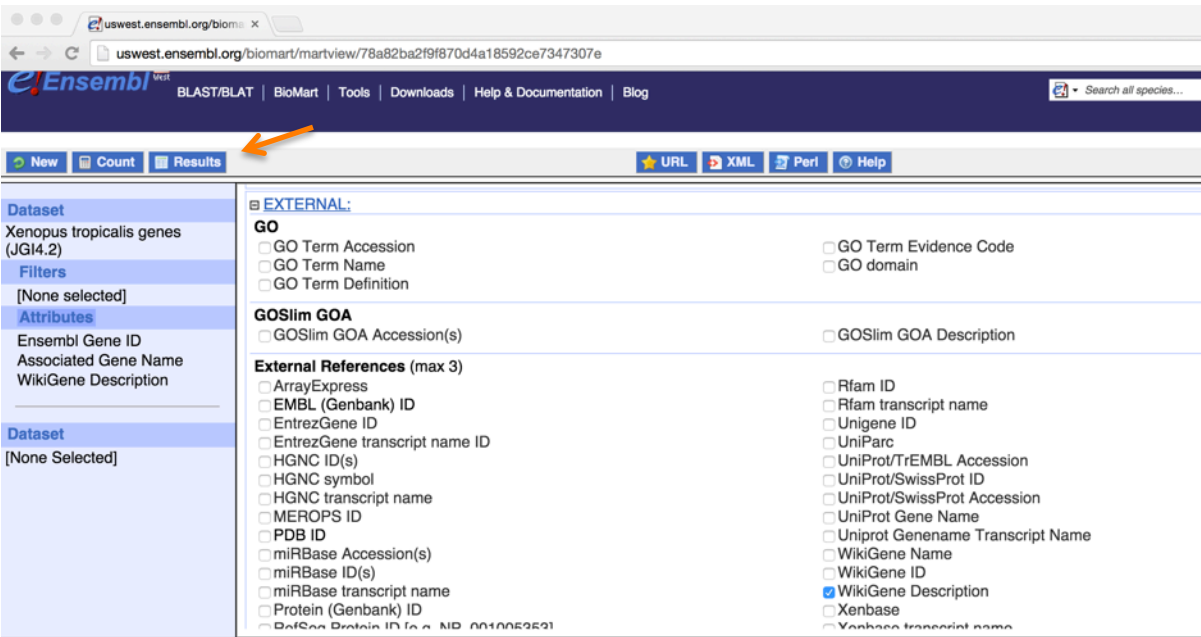


522

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



Step6. Click on “Results” icon.

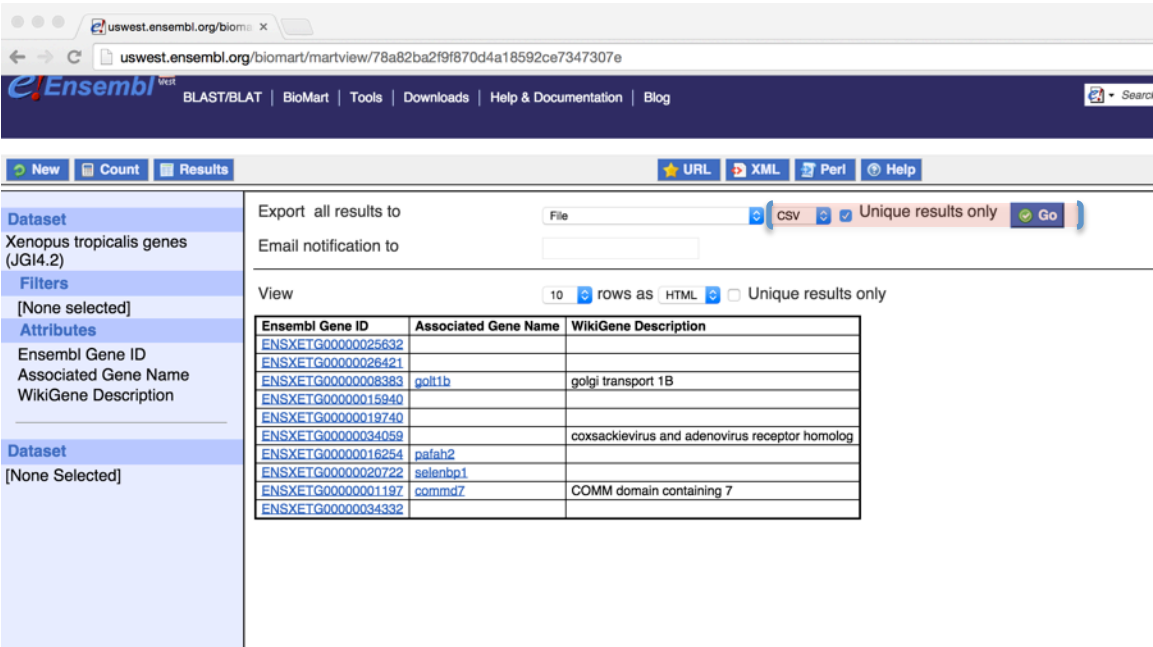


540

542

544

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



Step 8. Save and rename the result like “Xenopus.tropicalis\_gene\_name.txt”. There are three columns, separated by comma:

Ensembl Gene ID, Associated Gene Name, WikiGene Description  
 ENSXETG000000008383, golt1b, golgi transport 1B  
 ENSXETG000000034059, CARH, coxsackievirus and adenovirus receptor homolog  
 ENSXETG00000001197, commd7, COMM domain containing 7  
 .....

\*\*Repeat steps demonstrated above to download *Anolis carolinensis* reference protein from the Ensembl FTP and gene annotation file from the BioMart.

\*\*For this workshop, reference genomes and the corresponding biomart gene annotation files are already downloaded and located in “~/Desktop/MarkerDevelopment/associated\_data/”.

### First Annotation using *Xenopus tropicalis*:

#### Input:

1. A folder that contains trinity assemblies. This file “combined.fasta” is located in “~/Desktop/MarkerDevelopment/data/annotation\_pop/”

2. Reference protein downloaded from the ensemble:

Xenopus\_tropicalis.JGI\_4.2.pep.all.fa.

584

3. Reference biomart gene annotation file:

586

Xenopus\_tropicalis\_gene\_name.txt

588

4. Reference gtf file:

Anolis\_carolinensis.AnoCar2.0.78.gtf

590

### **Commands:**

592

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

594

ke@NGS:~/Desktop/MarkerDevelopment/data\$ 5-Annotation -a

596

~/Desktop/MarkerDevelopment/data/annotation\_pop/ -b

~/Desktop/MarkerDevelopment/associated\_data/Xenopus\_tropicalis.JGI\_4.2.pep.all.fa

598

-d ~/Desktop/SeqCap/programs/framedp-1.2.2/ -f

~/Desktop/MarkerDevelopment/associated\_data/Xenopus\_tropicalis\_gene\_name.txt -

600

n xenopus -e 1

602

#Repeat the same steps and command above to annotate assemblies using *Anolis carolinensis* as a reference

604

**##Copy the annotation results to “~/Desktop/MarkerDevelopment/data”**

606

ke@NGS:~/Desktop/MarkerDevelopment/data\$ scp -r

~/Desktop/MarkerDevelopment/associated\_data/annotation\_pop/\*

608

annotation\_pop/

610

### **Output:**

For each annotation, a new folder is generated under

612

“~/Desktop/MarkerDevelopment/data/annotation\_pop/”:

614

combined\_xenopus/

combined\_anole/

616

“combined” is the name of the assemblies and after that is the name of the reference used for annotate the assemblies.

618

620

##The annotated fasta files are “XXX\_xenopus\_annotated.fasta” and “XXX\_anole\_annotated.fasta”

622

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

624

combined\*/\*\_annotated.fasta

626

combined\_anole/combined\_anole\_annotated.fasta

combined\_xenopus/combined\_xenopus\_annotated.fasta

628

```

630  ##make a new folder "probe_design_pop" under
    "~/Desktop/MarkerDevelopment/data/".
632  ke@NGS:~/Desktop/MarkerDevelopment/data$ mkdir probe_design_pop

634  ##copy all the annotated fasta files to "probe_design_pop"
    ke@NGS:~/Desktop/MarkerDevelopment/data$ cp
636  annotation_pop/combined*/*annotated.fasta probe_design_pop/

638  ## check annotations in one of the annotated fasta files in post_annotation:
    ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design_pop$ head -4
640  combined_xenopus_annotated.fasta

642  >contig1    gs1_ge432    ENSXETG00000014175    vwa5a NA    5e-57
    TCTCTTACATGGACCCTTCC.....
644  >contig10    5u355_gs356_ge817_3u818 ENSXETG00000004176    mocs2
    molybdenum cofactor synthesis 2 2e-82
646  TGTGCACAGTGTGATGTAG.....

648  For contig1: "gs1" means coding region starts at position 1. "ge432" means coding
    region ends by position 432. No UTRs are present in this contig.
650  "ENSXETG00000014175" is the Ensembl gene ID obtained from Xenopus reference
    database. "vwa5a" is the gene name. "NA" is the wiki gene description and in this
652  case, wiki gene description is missing. "5e-57" is e-value in the BLAST search.

654  For contig10: "5u355" means 5UTR ends by position 355. "gs356" means coding
    region starts at position 356. "ge817" means coding region ends by position 817.
656  "3u818" means 3UTR starts at position 818. "ENSXETG00000004176" is the
    Ensembl gene ID obtained from Xenopus reference database. "mocs2" is the gene
658  name. "molybdenum cofactor synthesis 2" is the wiki gene description. "2e-82" is e-
    value in the BLAST search.
660
662  Run 5-Annotation with a GTF:

    Commands:
664  ke@NGS:~/Desktop/MarkerDevelopment$ 5-Annotation -a
    ~/Desktop/MarkerDevelopment/data/annotation_pop/ -b
666  ~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.pep.all.fa
    -d ~/Desktop/SeqCap/programs/framedp-1.2.2/ -n xenopus -g
668  ~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.78.gtf -e
    1
670
    The output by using GTF is slightly different since the header doesn't have gene
672  name descriptions. For example:
    >contig1    gs1_ge432    ENSXETG00000014175    vwa5a protein_coding    5e-
674  57 TCTCTTACATGGACCCTTCC.....

```

676 "gs1" means coding region starts at position 1. "ge432" means coding region ends by  
position 432. No UTRs are present in this contig. "ENSXETG00000014175" is the  
678 Ensembl gene ID obtained from Xenopus reference database. "vwa5a" is the gene  
name. **"protein\_coding" is the type of the gene.** "5e-57" is e-value in the BLAST  
680 search.

682

*\*6-ProcessAnnotation\*:*

684 Four sub-functions are included: Merge, Filter, TrimORF , Exon

686 If transcript targets are desired then run:  
Merge -> Filter -> TrimORF

688 If Exonic targets are desired then run:  
690 Merge -> Filter -> Exon

692 Now I will demonstrate pipelines for generating transcript targets (Merge -> Filter -  
> TrimORF).

694 “6-ProcessAnnotation Merge”: Merge annotations from various references. It also  
696 filters out redundant transcripts via self-blasting

698 Dependencies:  
BLAST+:  
700 [http://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE\\_TYPE=BlastDocs&DOC\\_TYPE=Download](http://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE_TYPE=BlastDocs&DOC_TYPE=Download)  
702 MUSCLE: <http://www.drive5.com/muscle/>  
cd-hit-est: <http://weizhongli-lab.org/cd-hit/>  
704

**Input:**

706 All annotated transcripts located in  
“~/Desktop/MarkerDevelopment/data/probe\_design\_pop/”

708

710 ##make a folder “probe\_design\_pop/”  
ke@NGS:~/Desktop/MarkerDevelopment/data\$ mkdir probe\_design\_pop/

712 ## copy data files to “probe\_design\_pop/”  
714 ke@NGS:~/Desktop/MarkerDevelopment/data\$ cp  
~/Desktop/MarkerDevelopment/associated\_data/annotation\_pop/combined\_\*/\*anno  
716 tated.fasta probe\_design\_pop/

718 ke@NGS:~/Desktop/MarkerDevelopment/data/probe\_design\_pop\$ ls  
combined\_anole\_annotated.fasta  
720 combined\_xenopus\_annotated.fasta

722 ##Make a new folder “other\_files” under  
“~/Desktop/MarkerDevelopment/data/probe\_design\_pop/”.

724 Use one of the annotated fasta files as a “primary” annotation file. Move the rest to a  
folder called “other\_files”. In the workshop we use  
726 “combined\_xenopus\_annotated.fasta” as the primary annotation file.

728 ke@NGS:~/Desktop/MarkerDevelopment/data/probe\_design\_pop\$ mkdir other\_files

```

730 ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design_pop$ mv
combined_anole_annotated.fasta other_files/
732
734 Commands:
# Run "6-ProcessAnnotation Merge":
736 ke@NGS:~/Desktop/MarkerDevelopment/data$ 6-ProcessAnnotation Merge -a
probe_design_pop/combined_xenopus_annotated.fasta -b
738 probe_design_pop/other_files/ -d frog
740 Output:
"1-frog_annotation_filtered.final" under
742 "~/Desktop/MarkerDevelopment/data/probe_design_pop/"
744
746 "6-ProcessAnnotation Filter": Basic filters on length, GC content, and repetitive
elements
748
Input:
750 "1-frog_annotation_filtered.final" produced by "6-ProcessAnnotation Merge"
752 Commands:
# Run "6-ProcessAnnotation Filter":
754 ke@NGS:~/Desktop/MarkerDevelopment/data$ 6-ProcessAnnotation Filter -f
probe_design_pop/1-frog_annotation_filtered.final -a 150 -b 1000 -R frogs -g frog
756
Output:
758 "2-frog_GC_length_repeatmasked.txt" under
~/Desktop/MarkerDevelopment/data/probe_design_pop/"
760
762 "6-ProcessAnnotation TrimORF": Trim off sequences outside the ORF
764 Input:
"2-frog_GC_length_repeatmasked.txt" produced by "6-ProcessAnnotation Filter":
766
Commands:
768 # Run "6-ProcessAnnotation TrimORF":
ke@NGS:~/Desktop/MarkerDevelopment/data$ 6-ProcessAnnotation TrimORF -f
770 probe_design_pop/2-frog_GC_length_repeatmasked.txt -b
~/Desktop/MarkerDevelopment/associated_data/combined_refProt.fasta -c 100 -o
772 frog
774

```

**Output:**

776 "3-frog\_transcript\_target.fasta" under  
 778 "~/Desktop/MarkerDevelopment/data/probe\_design\_pop/"

780 ++++++  
 782 Now I will demonstrate pipelines for generating exonic targets (Merge -> Filter -> Exon).

784 First run "6-ProcessAnnotation Merge" and "6-ProcessAnnotation Filter" following the instructions above

786  
 788 #Run 6-ProcessAnnotation Exon:

**Input:**

790 1. "2-frog\_GC\_length\_repeatmasked.txt" produced by "6-ProcessAnnotation Filter":

792 2. Combined protein reference of "Xenopus\_tropicalis.JGI\_4.2.pep.all.fa" and  
 794 "Anolis\_carolinensis.AnoCar2.0.pep.all.fa"

794 *ke@NGS:~/Desktop/MarkerDevelopment/associated\_data\$ cat*  
 796 *Anolis\_carolinensis.AnoCar2.0.pep.all.fa Xenopus\_tropicalis.JGI\_4.2.pep.all.fa >*  
*combined\_refProt.fasta*

798 3. Combined genomic DNA reference of  
 800 "Xenopus\_tropicalis.JGI\_4.2.dna\_rm.toplevel.fa" and  
 802 "Anolis\_carolinensis.AnoCar2.0.dna\_rm.toplevel.fa"  
*ke@NGS:~/Desktop/MarkerDevelopment/associated\_data\$ cat*  
*Xenopus\_tropicalis.JGI\_4.2.dna\_rm.toplevel.fa*  
*Anolis\_carolinensis.AnoCar2.0.dna\_rm.toplevel.fa > combined\_refGenome.fasta*

804

**Commands:**

806 # Run "6-ProcessAnnotation Exon" (do not execute the command!):  
 808 *ke@NGS:~/Desktop/MarkerDevelopment/data\$ 6-ProcessAnnotation Exon -p*  
*~/Desktop/MarkerDevelopment/associated\_data/combined\_refProt.fasta -g*  
 810 *~/Desktop/MarkerDevelopment/associated\_data/combined\_refGenome.fasta -f*  
*probe\_design\_pop/2-frog\_GC\_length\_repeatmasked.txt -e 150 -E 1000 -o frog*

812

**Output:**

814 "3-frog\_exon.fa"

816 #Let's copy the output file from  
 818 "~/Desktop/MarkerDevelopment/associated\_data/" to  
 "~/Desktop/MarkerDevelopment/data/probe\_design\_exons/".



```

820 ke@NGS:~/Desktop/MarkerDevelopment/data$ cp
    ~/Desktop/MarkerDevelopment/associated_data/3-frog_exon.fa probe_design_exons/
822
824 #check the output using "head"
826 >frog Contig1130_134_310      ENSXETG00000000011_exon2
    GGATCATGCCAAAGTTCTTCACTACATCGGAGCTGGGGTTGCCTTCCCAACCAGTATGTT
828 GTTCATTTTCTTTCAGTCTATCCTGACCTACCGCATGGCACACACTTATTGGAAGTGGTG
    GGCTGGACACGTACGCTGTCTTCTTACGTTGTTTGGACTGGTCATTTTAGTGCTTAG
830 >frog Contig2807_214_363      ENSXETG00000000013_exon11
    GATGCAAAAAATATAGAGGCAGAGGTAAGAGCTGAAGAAAGTTACAGATTTGAGCAT
832 TGTCCGGCTGCGGTTTACAGCATATTTGCCAGACAGCACTGGTGCCTATACTCTCCGTTT
    AAAGCCAGTCATTTCTGACCCTATCCATGAC
834 >frog Contig3806_208_390      ENSXETG00000000060_exon2
    GTCGTAGCCCTCCAGGTCTCTGAGCTTTTTTGATTTCAAAAAGGTTTCCCTCCACAGCCAG
836 CAGGGAGATTTGGAATCGCTGAGGATGCTCTGAGGTAACATACTGAGCTCCAGACAGT
    TCTCTTCCAGACGTAAACTTTGAGACGTGGACAATGTGAAACCTGAGCTGATATTTGG
838 GATAT
    .....
840
    "frog" is the name of the focal species
842 "Contig1130_134_310": This exon is located between position 134 to 310 of
    Contig1130.
844 "ENSXETG00000000011_exon2": Ensembl gene ID and numer of the exon

```