# 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 entire transcripts for marker development.  In this case please use "6-

48    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
     adapters, merging overlapping reads, removing duplicates and reads sourced from
140  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  *AATGATACGGCGACCACCGAGATCTACACcctgcgaACACTCTTTCCCTACACGACGCTCTTCCGATCT*
     *>P5_index2*
174  *AATGATACGGCGACCACCGAGATCTACACtgcagagACACTCTTTCCCTACACGACGCTCTTCCGATCT*
     *……*
176

     Note: The header of each adapter sequence has to be named strictly as "**P7_index**N"
178  or "**P5_index**N". 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

     *library       P7     P5*

186    *CGRL_index1    1*
       *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:
       *Escherichia coli* (bacteria + human + other genome resources if desired) genome in
194    fasta format.
       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
       "~/Desktop/MarkerDevelopment/data/":
202    *ke@NGS:~/Desktop/MarkerDevelopment/data$ mkdir cleaned_data*

204    #Run *2-ScrubReads*:
       *ke@NGS:~/Desktop/MarkerDevelopment/data$ 2-ScrubReads -f*
206    *~/Desktop/MarkerDevelopment/data/rawdata/raw/pre-clean/ -o*
       *~/Desktop/MarkerDevelopment/data/cleaned_data/ -a*
208    *~/Desktop/SeqCap/denovoTargetCapture/associated_files/Adapters.fasta -b*
       *~/Desktop/SeqCap/denovoTargetCapture/associated_files/libInfo.txt -t*
210    */home/ke/Desktop/SeqCap/programs/Trimmomatic-0.32/trimmomatic-0.32.jar -c*
       *~/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
       parameters when processing the real datasets.
216

       **Output:**
218    1. In "~/Desktop/MarkerDevelopment/data/cleaned_data/", six  .txt files per
       library are produced:
220     For example for library CGRL_index1, the six files are:
       CGRL_index1_1_final.txt (left reads)
222    CGRL_index1_2_final.txt (right reads)
       CGRL_index1_u_final.txt (merged or unpaired reads)
224    CGRL_index1.contam.out  (headers of reads aligned to bacteria)
       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
       find fastQC results for cleaned reads from each library.
230

*3-GenerateAssemblies*: Assemble RNAseq data using Trinity.

232

Dependencies:

234    Trinity http://trinityrnaseq.sourceforge.net

236    **Input:**
For each library, we will concatenate cleaned forward reads  (XXX_1_final.txt) and

238    unpaired reads (XXX_u_final.txt) and name the resulting read data file as
XXX_1_final.txt.

240

#Make a new folder called "raw_assembly" under

242    "~/Desktop/MarkerDevelopment/data/":
*ke@NGS:~/Desktop/MarkerDevelopment/data$ mkdir raw_assembly*

244

#Concatenate cleaned forward reads and unpaired reads and save them in

246    "raw_assembly":
*ke@NGS:~/Desktop/MarkerDevelopment/data$ cat*

248    *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*

250    *ke@NGS:~/Desktop/MarkerDevelopment/data$ cat*
*cleaned_data/CGRL_index15_1_final.txt cleaned_data/CGRL_index15_u_final.txt | sed*

252    *'s/\/2$/\/1/g' > raw_assembly/CGRL_index15_1_final.txt*
*ke@NGS:~/Desktop/MarkerDevelopment/data$ cat*

254    *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*

256

#Copy read2 of all libraries to "raw_assembly"

258    *ke@NGS:~/Desktop/MarkerDevelopment/data$ cp*
*cleaned_data/CGRL_index*_2_final.txt raw_assembly/*

260

**Commands:**

262    #Run Trinity on 4 processors.
 *ke@NGS:~/Desktop/MarkerDevelopment/data$ 3-GenerateAssemblies trinity -a*

264    *raw_assembly/ -c 5 -e 4*

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

268    **Output**:
There are quite a few intermediate files generated in

270    "~/Desktop/MarkerDevelopment/data/raw_assembly/CGRL_index1/".
"~/Desktop/MarkerDevelopment/data/raw_assembly/CGRL_index15/".

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

274

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

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

7

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 purpose of quick demonstration.  To better demonstrate how to use the next**
300 **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 examples 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:**
    #Display results in the output file "accuracy.out":
398 *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
assembled contig) and completeness (the percentage of expressed reference

408    transcripts covered by all matched assembled contigs)

410    **Commands:**
*ke@NGS:~/Desktop/MarkerDevelopment/data$ 4-AssemblyEvaluation  CONTIGUITY -*

412    *a annotation/ -b 300 -c*
*~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.cdna.all.f*

414    *a*

416    \*\*Note: that –b in function "CONTIGUITY" refers to the number of randomly selected
sequences from the reference protein database. In functions "BASIC",

418    "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":
*ke@NGS:~/Desktop/MarkerDevelopment/data/annotation$ less Contiguity.out*

424

    *Assemblies    complete(%)    contiguity(%)*

426    *CGRL_index1    13.46  11.47*
*CGRL_index15   23.36  23.36*

428    *GRL_index40   37.40  30.21*

430    _____

432    *5-Annotation*: annotate assembled contigs using a related reference protein
dataset that can be found in Ensembl Genome Browser

434    (http://www.ensembl.org/index.html)

436    Dependencies:
BLAST+:

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

440    FrameDP: https://iant.toulouse.inra.fr/FrameDP/cgi-bin/framedp.cgi?_wb_cfg=/www/iant/FrameDP/cgi-

442    bin/../cfg/FrameDP.cfg&_wb_session=WBuPAWHo&_wb_main_menu=Download&_wb_function=Download

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

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

448    modification of this script is needed.

450    ** Swiss-Prot (created in 1986) is a high quality manually annotated and non-redundant protein sequence database, which brings together experimental results,

452    computed features and scientific conclusions. UniProtKB/Swiss-Prot is now the
reviewed section of the UniProt Knowledgebase.

454

    **FrameDP: Sensitive peptide detection on noisy matured sequences. A self-training

456    integrative pipeline for predicting CDS in transcripts which can adapt itself to
different levels of sequence qualities.

458

    **Input:**

460    1. download a reference protein dataset from the Ensembl:

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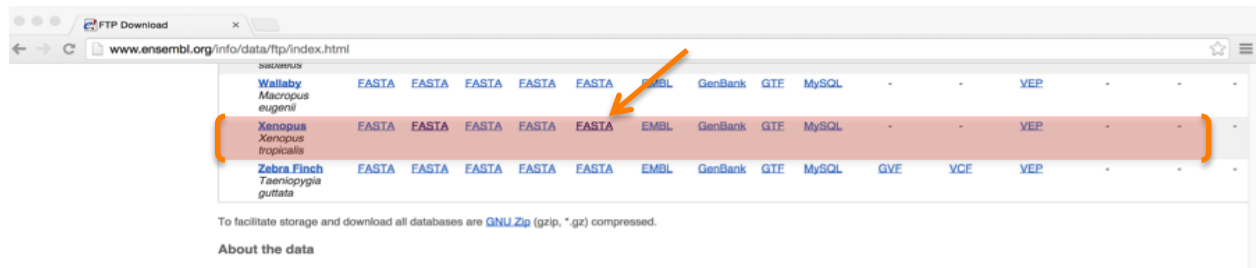"Download" located at the top.



464

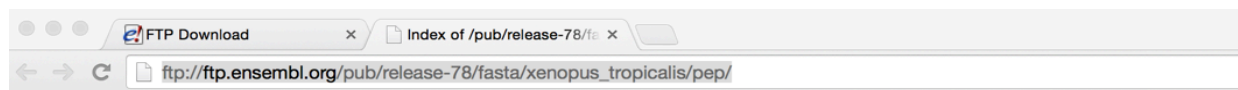Step2. Click on "Download data via FTP" to the left of the download page.

466



468

470     Step3. Select "All" in the "single species data" box in the FTP download page.



472

474

476

478

480

482

484

486

488

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



496

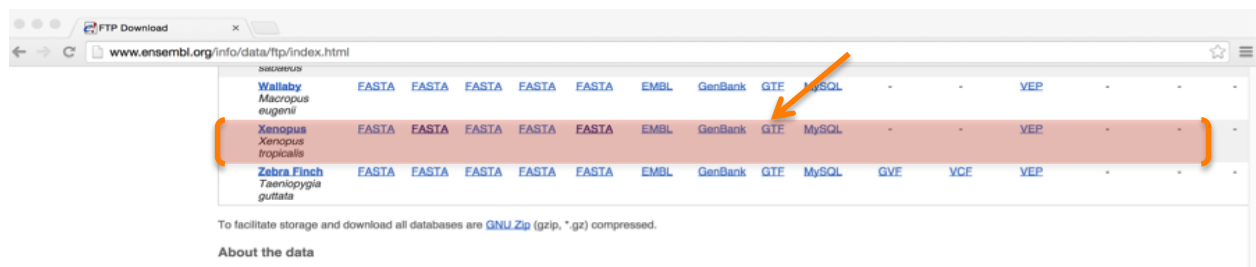498    Step 5: From FTP server, download reference protein fasta "XXX.pep.all. fa.gz"



504

506    Step 6: unzip the downloaded reference fasta: *gunzip*
       *Xenopus_tropicalis.JGI_4.2.pep.all. fa.gz*
508

       Step 7: Find and download the GTF (Gene transfer format (GTF) is a file format used
510    to hold information about gene structure) if there is one available for the reference.
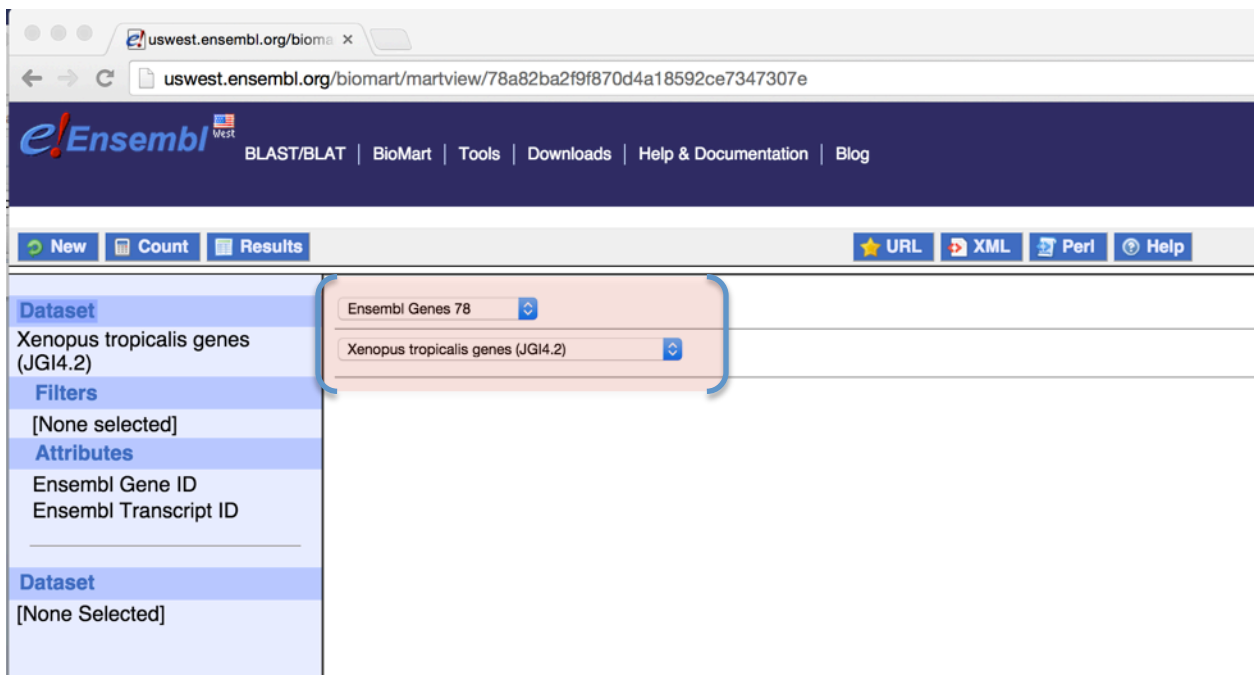       In this case we can see that *Xenopus tropicalis* has a GTF so we can download it.
512



514

516

       Step 8: unzip the downloaded GTF: *gunzip Xenopus_tropicalis.JGI_4.2.78.gtf.gz*
518

520    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
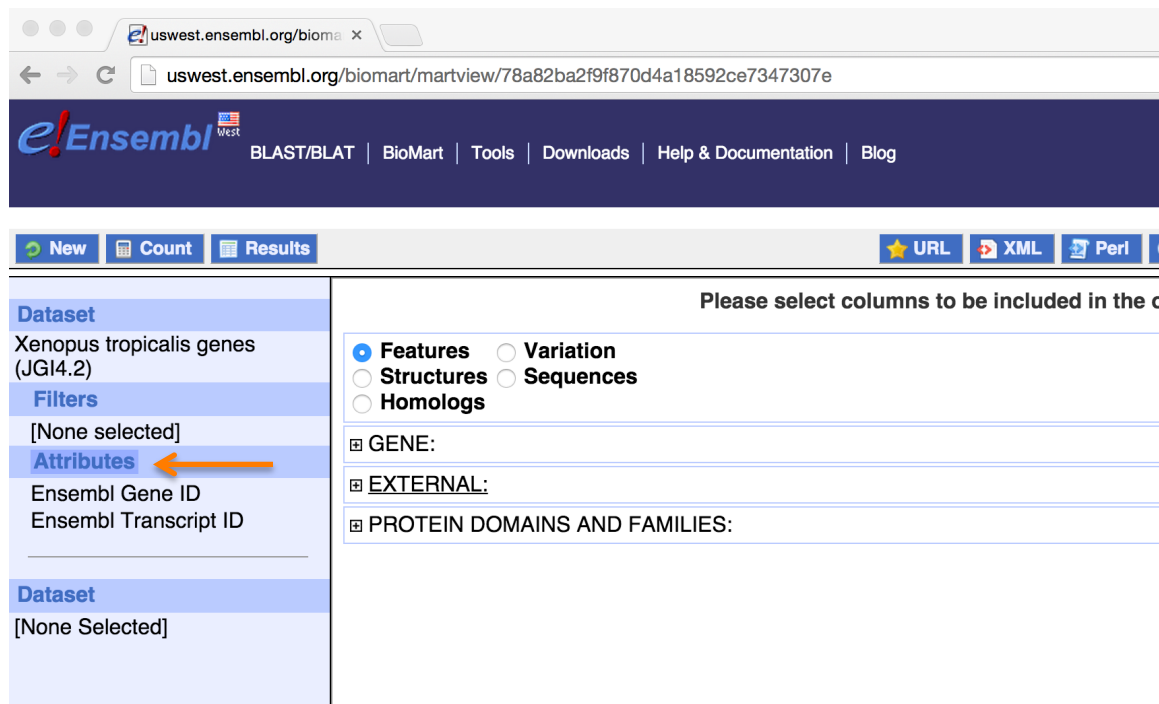522    file from the BioMart tool even though we have downloaded a GTF for the reference.

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



526

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



530

532

534

Step3.  Click on "Attributes" icon to the left.



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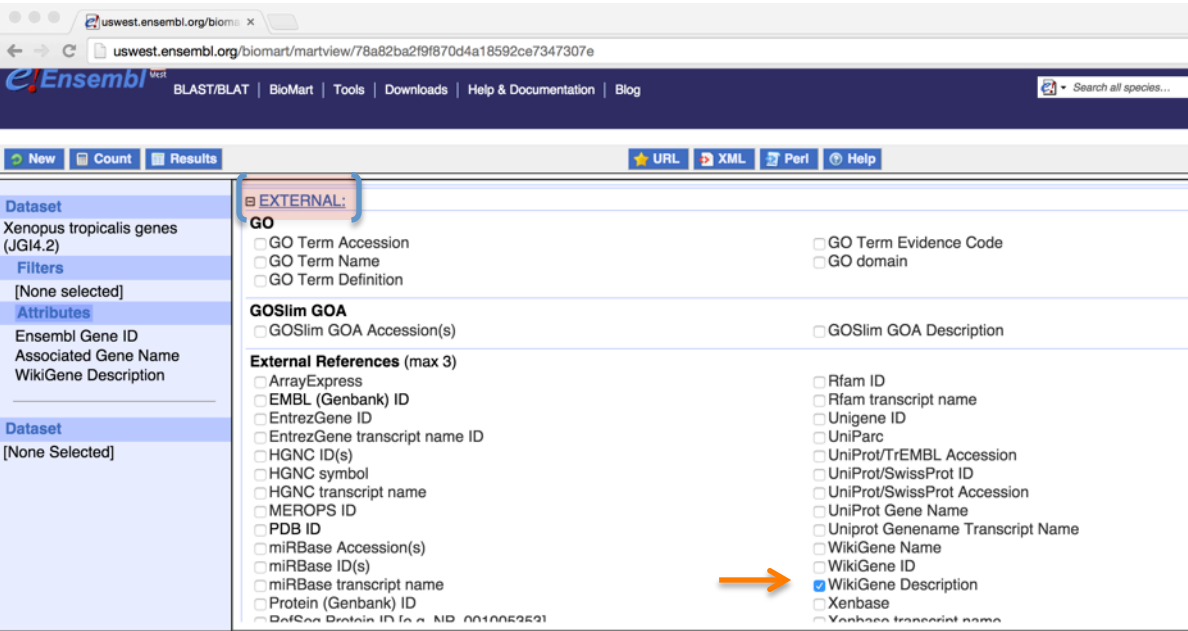


544

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"
572    box, and then click on "Go".



582

584    Step 8. Save and rename the result to be "*Xenopus.tropicalis* _gene_name.txt". There
       are three columns, separated by comma:
586
       Ensembl Gene ID, Associated Gene Name, WikiGene Description
588    ENSXETG00000008383, golt1b, golgi transport 1B
       ENSXETG00000034059, CARH, coxsackievirus and adenovirus receptor homolog
590    ENSXETG00000001197, commd7, COMM domain containing 7
       ......
592
       **For this workshop, a reference protein, a GTF and the corresponding biomart gene
594    name file are already downloaded and located in
       "~/Desktop/MarkerDevelopment/associated_data/".
596
       **Input:**
598    1. A folder that contains all trinity assemblies. These files are located in
       "~/Desktop/MarkerDevelopment/data/annotation/"
600
       2. Reference protein downloaded from the ensemble:
602    Xenopus_tropicalis.JGI_4.2.pep.all.fa.

604    3. Reference biomart gene annotation file:
       Xenopus_tropicalis_gene_name.txt
606
       *OR*

608    4. Reference GTF file:
*Xenopus_tropicalis.JGI_4.2.78.gtf*

610

**Commands:**

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

614

*ke@NGS:~/Desktop/MarkerDevelopment/data$ 5-Annotation  -a*
616    *~/Desktop/MarkerDevelopment/data/annotation/ -b*
*~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.pep.all.fa*
618    *-d ~/Desktop/SeqCap/programs/framedp-1.2.2/ -f*
*~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis_gene_name.txt -*
620    *n xenopus -e 1*

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

626    **Output:**
For each individual trinity assembly, a new folder is generated under
628    "~/Desktop/MarkerDevelopment/data/annotation/":

630    CGRL_index1_xenopus/
CGRL_index14_xenopus /
632    CGRL_index40_xenopus /

634    ##The annotated fasta files are named as "XXX_xenopus_annotated.fasta".

636    *ke@NGS:~/Desktop/MarkerDevelopment/data/annotation$ ls*
*CGRL_index*/*annotated.fasta*
638

*CGRL_index15_xenopus/CGRL_index15_xenopus_annotated.fasta*
640    *CGRL_index50_xenopus/CGRL_index50_xenopus_annotated.fasta*
*CGRL_index1_xenopus/CGRL_index1_xenopus_annotated.fasta*
642

##make a new folder "probe_design" under
644    "~/Desktop/MarkerDevelopment/data/".
*ke@NGS:~/Desktop/MarkerDevelopment/data$ mkdir probe_design*
646

##copy all the annotated fasta files to *"probe_design"*
648    *ke@NGS:~/Desktop/MarkerDevelopment/data$ cp*
*annotation/CGRL_index*/*annotated.fasta  probe_design/*
650

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

654    >contig1    gs1_ge432    ENSXETG00000014175    vwa5a NA    5e-57
       TCTCTTACATGGACCCTTCC......
656    >contig10    5u355_gs356_ge817_3u818 ENSXETG00000004176    mocs2
       molybdenum cofactor synthesis 2 2e-82
658    TGTGCACAGTGTGATGTAG......

660    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.
662    "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
664    case, wiki gene description is missing.  "5e-57" is e-value in the BLAST search.

666    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.
668    "3u818" means 3UTR starts at position 818.  "ENSXETG00000004176" is the
       Ensembl gene ID obtained from Xenopus reference database.  "mocs2" is the gene
670    name. "molybdenum cofactor synthesis 2" is the wiki gene description.  "2e-82" is e-
       value in the BLAST search.
672
       ~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~
674    Run 5-Annotation with a GTF

676    **Commands:**
       *ke@NGS:~/Desktop/MarkerDevelopment$ 5-Annotation -a*
678    *~/Desktop/MarkerDevelopment/data/annotation/ -b*
       *~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.pep.all.fa*
680    *-d ~/Desktop/SeqCap/programs/framedp-1.2.2/  -n xenopus -g*
       *~/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.78.gtf -e*
682    *1*

684    The output by using GTF is slightly different since the header doesn't have gene
       name descriptions. For example:
686    >contig1    gs1_ge432    ENSXETG00000014175    vwa5a protein_coding    5e-
       57
688    TCTCTTACATGGACCCTTCC......

690    "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
692    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
694    search.

696

698

20

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:
706   BLAST+:
http://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE_TYPE=BlastDocs&DOC_TYPE=Downlo
708   ad
MUSCLE: http://www.drive5.com/muscle/
710   cd-hit-est: http://weizhongli-lab.org/cd-hit/

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

716   **Input:**
All annotated transcripts located in
718   "~/Desktop/MarkerDevelopment/data/probe_design"

720   ##
*ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design$ ls*
722   *CGRL_index15_xenopus_annotated.fasta*
*CGRL_index40_xenopus_annotated.fasta*
724   *CGRL_index1_xenopus_annotated.fasta*

726   Make a new folder "other_files" under
"~/Desktop/MarkerDevelopment/data/probe_design/".
728   Use one of the annotated files as a "primary" annotation file. Move the rest to a
folder "other_files". In the workshop we use CGRL_index1 as the "primary"
730   annotation file.

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

734   *ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design$ mv*
*CGRL_index15_xenopus_annotated.fasta  CGRL_index40_xenopus_annotated.fasta*
736   *other_files/*

738   **Commands:**
# 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*
742   *1000*

744

**Output:**

746      #Under "ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design/" a new folder called "results" was created by the script.

748      *ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design$ cd results/*

750      #Markers that passed all filters are stored in "marker_kept.txt".  First take a How many markers are kept?

752      *ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design/results$ wc -l marker_kept.txt*

754      *1050 marker_kept.txt*

756      *ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design/results$ less -S marker_kept.txt*

758

     **Transcript_name**: Ensembl Gene ID

760      **avgDiv**: Average sequence divergence (avg. %mismatches)
     **varianceDiv**: Variance of sequence divergence

762      **avgLength**: Average length of the marker
     **avgGC**: Average CG content of the marker

764      **div_CGRL_index15_xenopus_annotated _vs_CGRL_index1_xenopus_annotated**: sequence divergence  between CGRL_index15 and CGRL_index1

766      **div_CGRL_index15_xenopus_annotated _vs_CGRL_index40_xenopus_annotated**: sequence divergence  between CGRL_index15 and CGRL_index40

768      **div_CGRL_index1_xenopus_annotated _vs_CGRL_index40_xenopus_annotated**: sequence divergence  between CGRL_index1 and CGRL_index40

770

772      #Select the markers that you would like to use for probe design. In this case choose 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*
776      *marker_kept.txt > marker_final.txt*

778      +++++++++++++++++++++++++++++++++
     Now we will use command "6-MarkerSelectionTRANS seq" to generate input fasta
780      files for probe design:

782      **Input:**
     1.  A final set of markers you would like to use for probe design -
784      "~/Desktop/MarkerDevelopment/data/probe_design/results/marker_final.txt"

786      2. A folder containing all trimmed transcripts in fasta format. These files were created by "6-MarkerSelectionTRANS markers" and are named as XXX.final2 –
788      "~/Desktop/MarkerDevelopment/data/probe_design/results/"

790

     **Commands:**

792      # Run *6-MarkerSelectionTRANS seq:*
     *ke@NGS:~/Desktop/MarkerDevelopment/data$ 6-MarkerSelectionTRANS seq -f*
794      *probe_design/results/marker_final.txt -d probe_design/results/*
     *The target size for CGRL_index15_xenopus _annotated.final2 is 700532bp!*
796      *The target size for CGRL_index1_xenopus _annotated.final2 is 701541bp!*
     *The target size for CGRL_index40_xenopus _annotated.final2 is 701593bp!*
798

     **Output:**
800      #A new folder "Probe_Design" was created by the script. cd to this folder:
     *ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design/results$ cd*
802      *Probe_Design/*

804      #Three fasta sequence files contain sequences of orthologous markers are
     generated and ready for submission for probe design:
806      *ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design/results/Probe_Design$*
     *ls *exonic_targets.txt*
808

     *CGRL_index15_xenopus _annotated_exonic_targets.txt*
810      *CGRL_index1_xenopus _annotated_exonic_targets.txt*
     *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
cd-hit-est: http://weizhongli-lab.org/cd-hit/

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

824     ad
MUSCLE: http://www.drive5.com/muscle/

826

We will run "6-MarkerSelectionEXONS exons" to identify orthologous exons in

828     transcriptomes from each of the species.

830     **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

832     reference to identify ortholgous exons from each of the transcriptomes.

834     **However, if a .gtf file is available then I recommend first run "ParseGTF" to obtain exonic sequences from the reference and then run "6-MarkerSelectionEXONS

836     exons".

838

++++++++++++++++++++++++++++++++++++++++++++++++++++++++

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

842

**Input:**

844     1. Under "~/Desktop/MarkerDevelopment/data/" make a new folder "probe_design_exons/":

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

848     2. copy all annotated transcripts to "~/Desktop/MarkerDevelopment/data/probe_design_exons".

850

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

852     *probe_design/CGRL_index1_xenopus_annotated.fasta*
*probe_design/other_files/CGRL_index\* probe_design_exons/*

854

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

856

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

858     *CGRL_index15 _xenopus _annotated.fasta*
*CGRL_index40 _xenopus _annotated.fasta*

860     *CGRL_index1 _xenopus _annotated.fasta*

862

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

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

868 Both 3 and 4 can be downloaded through Ensembl following the instruction above.
For this workshop these two files are located under
870 "~/Desktop/MarkerDevelopment/associated_data".

872 **Command:**
#Run "6-MarkerSelectionEXONS exons"
874 *ke@NGS:~/Desktop/MarkerDevelopment/data$ 6-MarkerSelectionEXONS exons -p
'/home/ke/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.*
876 *pep.all.fa' -g
'/home/ke/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.*
878 *dna_rm.nonchromosomal.fa' -f
~/Desktop/MarkerDevelopment/data/probe_design_exons -E 1000*
880
** *"6-MarkerSelectionEXONS exons"* takes very long time to run so please do not run
882 it during the workshop.  Let's skip this step and copy the output files directly from
"associated_data":
884
*ke@NGS:~/Desktop/MarkerDevelopment/data$ cp
886 ~/Desktop/MarkerDevelopment/associated_data/probe_design_exons/*.nr
~/Desktop/MarkerDevelopment/associated_data/probe_design_exons/marker_*
888 probe_design_exons/*

890

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

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

900 **exon_name**: Exon ID
**avgDiv**: Average sequence divergence (avg. %mismatches)
902 **varianceDiv**: Variance of sequence divergence
**avgLength**: Average length of the exons
904 **avgGC**: Average CG content of the exons
**div_CGRL_index15_xenopus _annotated_vs_CGRL_index1_xenopus _annotated**:
906 sequence divergence  between CGRL_index15 and CGRL_index1

908 **div_CGRL_index15_xenopus _annotated_vs_CGRL_index40_xenopus _annotated**: sequence divergence  between CGRL_index15 and CGRL_index40

**div_CGRL_index1_xenopus _annotated_vs_CGRL_index40_xenopus _annotated**:
910 sequence divergence  between CGRL_index1 and CGRL_index40


912


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


918

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


926 **do not run ParseGTF in the workshop


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*
948 *probe_design/CGRL_index1_xenopus_annotated.fasta*
*probe_design/other_files/CGRL_index* probe_design_exons/*
950


952 **Command:**

#run "6-MarkerSelectionEXONS exons" (do not run it in the workshop)

954 *ke@NGS:~/Desktop/MarkerDevelopment/data$ 6-MarkerSelectionEXONS exons -p '/home/ke/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.*

956 *pep.all.fa' -g '/home/ke/Desktop/MarkerDevelopment/associated_data/Xenopus_tropicalis.JGI_4.2.*

958 *dna_rm.toplevel.fa' -f ~/Desktop/MarkerDevelopment/data/probe_design_exons -E 1000*

960

**Output:**

962 Same as above:

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

966

**Input:**

968 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 -

970 "~/Desktop/MarkerDevelopment/data/probe_design/results/ marker_kept_one_exon_per_gene.txt"

972

2. A folder containing non-redundant exonic markers in fasta format. These files

974 were created by "6-MarkerSelectionEXONS exons" and are named as XXX _exon.fa.nr – "~/Desktop/MarkerDevelopment/data/probe_design_exons/"

976

**Commands:**

978 # Run *6-MarkerSelectionEXONS seq:*

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

982

**Output:**

984 #A new folder "Probe_Design" was created by the script.  cd to this folder: *ke@NGS:~/Desktop/MarkerDevelopment/data/probe_design_exons$ cd*

986 *Probe_Design/*

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

990

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

992 *exonic_targets.txt*

994 *CGRL_index15_xenopus _annotated_exonic_targets.txt*
*CGRL_index1_xenopus _annotated_exonic_targets.txt*

996 *CGRL_index40_xenopus _annotated_exonic_targets.txt*