# Pipeline for *de novo* Targeted Capture

October 6, 2014

Contributors: Sonal Singhal, Ke Bi, Tyler Linderoth

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/MVZSEQ/denovoTargetCapture
_____

Scripts included in this pipeline:

1-PreCleanup

2-ScrubReads

3-GenerateAssemblies

4-FinalAssembly

5-FindTargets

6-AssemblyEvaluation (optional)

7-Alignment

8-ExonCaptureEvaluation (optional)

9-preFiltering

10-SNPcleaner

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

- Exon/exome/sequence capture dataset
   Use 1->2->3->4->5->6->7->8 when no reference genome is available;
   Use 1->2->7->8 when a reference genome is available

50    - Genomic dataset (with a pre-existing genome or *de novo* genome scaffolds)
          Use 1->2->7
52

        - *de novo* transcriptome dataset
54        Use 1->2->3->4->5->7

56    - Single RAD/GBS dataset
          Use 1->2->3->RAD/GBS Tag filtering->7
58

        - ddRAD/ddGBS dataset
60        Use 1->2->ddRAD pipelines->7

62    - UCE dataset
          Use 1->2->UCE pipelines->7
64    _____

66    *1-PreCleanup*: Reformats raw sequencing reads from Illumina HiSeq or MiSeq for
        *2-ScrubReads*. Specifically, in this step we will remove reads that did not pass the
68    Illumina quality control filters and modify the sequence identifiers.

70    Dependencies:
        FastQC: http://www.bioinformatics.babraham.ac.uk/projects/fastqc/
72

        **Input:**
74    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
76    saved under "/home/ke/Desktop/SeqCap/data/rawdata/library/". Compressed
        fastq sequence files are saved in each of these folders.
78

        Fastq files use the following naming scheme:
80    <sample name>_<barcode sequence>_L<lane (0-padded to 3 digits)>_R<read
        number>_<set number (0-padded to 3 digits)>.fastq.gz
82

        For example, in "CGRL_index15_CGACCTG_L006_R1_001.fastq.gz":

84    sample name:  CGRL_index15
        barcode sequence:  CGACCTG
86    lane (0-padded to 3 digits): 006
        read number: 1
88    set number (0-padded to 3 digits):  001

90    #Make a new folder called "raw" under "~/Desktop/SeqCap/data/rawdata/":
        *ke@NGS:~/Desktop/SeqCap/data/rawdata$ mkdir raw*
92

        #Copy all these compressed fastq files from each folder (CGRL_index1,

94    CGRL_index15, CGRL_index40) to "raw":
*ke@NGS:~/Desktop/SeqCap/data/rawdata$ cp library/CGRL_index*/*.gz raw/*

96

#Check data files in "raw":
98    *ke@NGS:~/Desktop/SeqCap/data/rawdata$ ls raw/**
*CGRL_index15_CGACCTG_L006_R1_001.fastq.gz*
100   *CGRL_index15_CGACCTG_L006_R2_001.fastq.gz*
*CGRL_index1_TCGCAGG_L006_R1_001.fastq.gz*
102   *CGRL_index1_TCGCAGG_L006_R2_001.fastq.gz*
*CGRL_index40_TTCGCAA_L006_R1_001.fastq.gz*
104   *CGRL_index40_TTCGCAA_L006_R2_001.fastq.gz*


106

**Commands:**
108   #cd to the working directory:
*ke@NGS:~/Desktop/SeqCap/data/rawdata$ cd ..*
110

#run  1-PreCleanup with fastq evaluation
112   *ke@NGS:~/Desktop/SeqCap/data$ 1-PreCleanup*
*~/Desktop/SeqCap/data/rawdata/raw/ fastqc*
114

**Output:**
116   Three new folders will be created under "~/Desktop/SeqCap/data/rawdata/raw/":
"pre-clean"
118   "combined"
"pre-clean/evaluation"
120

- Folder "pre-clean" contains reformatted raw fastq reads.
122   CGRL_index1_R1.fq
CGRL_index1_R2.fq
124   CGRL_index15_R1.fq
CGRL_index15_R2.fq
126   CGRL_index40_R1.fq
CGRL_index40_R2.fq
128

- Folder "combined" contains merged, compressed, fastq data files (not used by the
130   following pipeline).
CGRL_index1_TCGCAGG_L006_R1.fastq.gz
132   CGRL_index1_TCGCAGG_L006_R2.fastq.gz
CGRL_index15_CGACCTG_L006_R1.fastq.gz
134   CGRL_index15_CGACCTG_L006_R2.fastq.gz
CGRL_index40_TTCGCAA_L006_R1.fastq.gz
136   CGRL_index40_TTCGCAA_L006_R2.fastq.gz

138   - Folder "evaluation" contains fastQC results for each data file.
CGRL_index1_R1.fq_fastqc/

140  CGRL_index1_R2.fq_fastqc/
     CGRL_index15_R1.fq_fastqc/
142  CGRL_index15_R2.fq_fastqc/
     CGRL_index40_R1.fq_fastqc/
144  CGRL_index40_R2.fq_fastqc/

146  Questions:
     1. Check the sequence identifiers and the number of reads in fastq files before and
148  after running *1-PreCleanup* and compare the results.
     2. Check the fastQC evaluation results for the raw data

150

     _____

152
     *2-ScrubReads*: Clean up raw data, which includes trimming for quality, removing
154  adapters, merging overlapping reads, removing duplicates and reads sourced from
     contamination

156
     Dependencies:
158  cutadapt: http://code.google.com/p/cutadapt/
     COPE: http://sourceforge.net/projects/coperead/
160  Bowtie2: http://sourceforge.net/projects/bowtie-bio/files/bowtie2/
     FastQC: http://www.bioinformatics.babraham.ac.uk/projects/fastqc/
162  FLASh-modified: modified version of FLASh by Filipe G. Vieira.
     https://github.com/MVZSEQ/Exon-capture
164  Trimmomatic: http://www.usadellab.org/cms/?page=trimmomatic

166  **Input:**
     1. Reformatted fastq files created by *1-PreCleanup*:
168  #Check the raw data files:
     *ke@NGS:~/Desktop/SeqCap/data/rawdata/raw/pre-clean$ ls *.fq*
170  *CGRL_index1_R1.fq*
     *CGRL_index1_R2.fq*
172  *CGRL_index15_R1.fq*
     *CGRL_index15_R2.fq*
174  *CGRL_index40_R1.fq*
     *CGRL_index40_R2.fq*
176
     2. A fasta file that contains adapter sequences:
178  #Check the format of adapter sequence file:
     *ke@NGS:~/Desktop/SeqCap/denovoTargetCapture/associated_files $ less -S*
180  *Adapters.fasta*
     *>P7_index1*
182  *CAAGCAGAAGACGGCATACGAGATcctgcgaGTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT*
     *>P7_index2*
184  *CAAGCAGAAGACGGCATACGAGATtgcagagGTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT*
     *......*
186  *>P5_index1*

4

188 *AATGATACGGCGACCACCGAGATCTACACcctgcgaACACTCTTTCCCTACACGACGCTCTTCCGATCT*
    *>P5_index2*
    *AATGATACGGCGACCACCGAGATCTACACtgcagagACACTCTTTCCCTACACGACGCTCTTCCGATCT*
190 *......*

192 Note: The header of each adapter sequence has to be named strictly as "**P7_index**N"
    or "**P5_index**N". N is the number of index. It is OK to put all adapters in this file but
194 your libraries only use a subset of them.

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

200 *library        P7      P5*
    *CGRL_index1    1*
202 *CGRL_index15   15*
    *CGRL_index40   40*
204

    leave the "P5" column blank if you only have indexes in P7 adapters in the libraries.
206

    4. Contaminant file:
208 *Escherichia coli* ( + human + other genome resources if desired) genome in fasta
    format.
210 This file (e_coli_K12.fasta) is saved in
    "~/Desktop/SeqCap/denovoTargetCapture/associated_files/ecoli/"
212


214 **Commands:**
    #Make a new folder called "cleaned_data" in "~/Desktop/SeqCap/data/":
216 *ke@NGS:~/Desktop/SeqCap/data$ mkdir cleaned_data*

218 #Run *2-ScrubReads*:
    *ke@NGS:~/Desktop/SeqCap/data$ 2-ScrubReads -f*
220 *~/Desktop/SeqCap/data/rawdata/raw/pre-clean/ -o*
    *~/Desktop/SeqCap/data/cleaned_data/ -a*
222 *~/Desktop/SeqCap/denovoTargetCapture/associated_files/Adapters.fasta -b*
    *~/Desktop/SeqCap/denovoTargetCapture/associated_files/libInfo.txt -t*
224 */home/ke/Desktop/SeqCap/programs/Trimmomatic-0.32/trimmomatic-0.32.jar -c*
    *~/Desktop/SeqCap/denovoTargetCapture/associated_files/ecoli/e_coli_K12.fasta -e*
226 *200 -m 15 -z*

228 Note: I use default values for most of the arguments. Users should adjust these
    parameters when processing the real datasets.
230

    **Output:**
232 1. In "~/Desktop/SeqCap/data/cleaned_data/", six  .txt files per library are
    produced:

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

242    2. In "~/Desktop/SeqCap/data/cleaned_data/evaluation/", you can find fastQC
       results for cleaned reads from each library.
244
       Questions:
246    1. Check the %reads that are exact duplicates, %reads that are likely derived from
       microbial genome and %reads that contain low complexity.
248    2. Check the fastQC evaluation results of cleaned reads and then compare them to
       those of raw reads. Is the quality improved?
250    _____

252    *3-GenerateAssemblies*: Assemble sequence capture data using ABySS.

254    We use a multiple-kmer approach to assemble our data. If there is even coverage
       and even polymorphism levels across the assembled genome, there should (in
256    theory) be one k-mer that best assembles the data. In reality, coverage and
       polymorphism vary across captured loci, and using multiple k-mers is a way to bet
258    hedge and get good assemblies for all loci. In assembling your data, it is important to
       consider which samples to use in your assembly. Ideally, you could assemble across
260    multiple individuals to increase your read depth, and thus, assembly contiguity and
       continuity. However, for many projects, more individuals can also mean increased
262    polymorphism. While we have found the assemblers are more robust to
       polymorphism than the program writers themselves often suggest, increased
264    polymorphism does lead to shorter contigs and increased misassemblies. With these
       sample data, we assembled across all in-group samples – this seemed like the best
266    balance between having enough data to power assembly while not introducing too
       much polymorphism.
268
       Dependencies:
270    ABySS (compiled with OpenMPI and Google sparsehash):
       http://www.bcgsc.ca/platform/bioinfo/software/abyss
272
       **Input:**
274    Concatenated cleaned reads from libraries that you would like to assemble together.
       The libraries to be assembled together have to be genetically similar: ideally,
276    samples from the same population.  In this example we want to assemble
       CGRL_index1, CGRL_index15 and CGRL_index40 together.
278
       #Make a new folder called "raw_assembly" under "~/Desktop/SeqCap/data/":

280 *ke@NGS:~/Desktop/SeqCap/data$ mkdir raw_assembly*

282 #Concatenate cleaned reads and save them in "raw_assembly":
 *ke@NGS:~/Desktop/SeqCap/data$ cat cleaned_data/CGRL_index*_1_final.txt >*
284 *raw_assembly/combined_1_final.txt*
 *ke@NGS:~/Desktop/SeqCap/data$ cat cleaned_data/CGRL_index*_2_final.txt >*
286 *raw_assembly/combined_2_final.txt*
 *ke@NGS:~/Desktop/SeqCap/data$ cat cleaned_data/CGRL_index*_u_final.txt >*
288 *raw_assembly/combined_u_final.txt*

290 #Inside "raw_assemblies" make a new folder "results":
 *ke@NGS:~/Desktop/SeqCap/data$ mkdir raw_assembly/results*
292

294 **Commands:**
 #Run ABySS on two processors using kmer sizes of 21, 31, 41, 51, 61, and 71.
296 *ke@NGS:~/Desktop/SeqCap/data$ 3-GenerateAssemblies abyss -reads*
 *~/Desktop/SeqCap/data/raw_assembly/ -mpi /usr/bin/mpirun -out*
298 *~/Desktop/SeqCap/data/raw_assembly/results/ -kmer 21 31 41 51 61 71 -np 2*

300 Note: Your labtop will not be able to handle memory intensive ABySS assemblies.

302 **Output**:
 There are a lot of intermediate files created in
304 "~/Desktop/SeqCap/data/raw_assembly/results/combined/".

306 #To show the assemblies that we need for the next step:
 *ke@NGS:~/Desktop/SeqCap/data$ cd raw_assembly/results/combined/*
308 *ke@NGS:~/Desktop/SeqCap/data/raw_assembly/results/combined$ ls *-contigs.fa*
 *combined_k21_cov_default-contigs.fa*
310 *combined_k31_cov_default-contigs.fa*
 *combined_k41_cov_default-contigs.fa*
312 *combined_k51_cov_default-contigs.fa*
 *combined_k61_cov_default-contigs.fa*
314 *combined_k71_cov_default-contigs.fa*

316 #Combine all the raw assemblies and write the result to a new file called
 "all_assemblies.fasta":
318 *ke@NGS:~/Desktop/SeqCap/data/raw_assembly/results/combined$ cat*
 *combined_*_cov_default-contigs.fa > all_assemblies.fasta*
320

 #Make a new folder called "merge_assemblies" under "~/Desktop/SeqCap/data/":
322 *ke@NGS:~/Desktop/SeqCap/data $ mkdir merge_assemblies*

324 #Copy "all_assemblies.fasta" into "merge_assemblies/":

326 *ke@NGS:~/Desktop/SeqCap/data $ cp raw_assembly/results/combined/all_assemblies.fasta merge_assemblies*

328 _____

330 *4-FinalAssembly*: Combining assembled contigs across multiple k-mers to generate a final assembly introduces a lot of redundancy into the final assembly. To address
332 this, we use a lightweight assembler cap3 and other programs (blat, cd-hit-est) to merge contigs and to remove redundancies.
334

Dependencies:
336 CAP3: http://seq.cs.iastate.edu/cap3.html
blat: http://users.soe.ucsc.edu/~kent/src/
338 cd-hit-est: https://code.google.com/p/cdhit/downloads/list

340 **Input:**
Concatenated raw assemblies
342 "~/Desktop/SeqCap/data/merge_assemblies/all_assemblies.fasta" produced by *3-GenerateAssemblies*
344

**Commands:**
346 *ke@NGS:~/Desktop/SeqCap/data$ 4-FinalAssembly -a ~/Desktop/SeqCap/data/merge_assemblies/ -c 1000*
348

Note: when analyzing real data, users should test these parameters (-d -e -b) for
350 optimal results.

352 **Output:**
Several files are created in "~/Desktop/SeqCap/data/merge_assemblies/".  The
354 data file that we need for the next step is "all_assemblies.fasta.final".

356 # Rename "all_assemblies.fasta.final":
*ke@NGS:~/Desktop/SeqCap/data/merge_assemblies$ mv all_assemblies.fasta.final*
358 *all_assemblies_final.fasta*

360 _____

362 *5-FindTargets*: identify contigs that are stemmed from the targeted loci and use these contigs as a reference (aka. a pseudo-reference)
364

Here, we suggest taking a very conservative approach to define the reference
366 genome against which you will align your reads. You will likely get many multiples more contigs than loci you targeted. Some of these might be junk; some might be
368 real. Rather than try to identify which of the extraneous contigs are junk or real, we suggest using only those contigs which match to the original targets. To do so, we
370 implement a BLAST approach, which identifies which contig has the best-hit match

to one's targeted loci.

372

Dependencies:

374 blastn (BLAST+): ftp://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/LATEST/
cd-hit-est: https://code.google.com/p/cdhit/downloads/list

376

**Input:**

378 1. "~/Desktop/SeqCap/data/merge_assemblies/all_assemblies_final.fasta"
produced by *4-FinalAssembly*.

380

2. Targeted loci fasta file:

382 "~/Desktop/SeqCap/denovoTargetCapture/original_target /targeted_loci.fasta"
contains loci/genes/exons from which probes are designed.

384

**Commands:**

386 #Make a new folder called "in_target_assemblies": under
"~/Desktop/SeqCap/data/":

388 *ke@NGS:~/Desktop/SeqCap/data$ mkdir in_target_assemblies*

390 #Run *5-FindingTargets* in "~/Desktop/SeqCap/data/in_target_assemblies/":
*ke@NGS:~/Desktop/SeqCap/data/in_target_assemblies$ 5-FindingTargets -t*

392 *~/Desktop/SeqCap/denovoTargetCapture/original_target/targeted_loci.fasta -a*
*~/Desktop/SeqCap/data/merge_assemblies/all_assemblies_final.fasta -o*

394 *in_target.fasta -e in_target.captured*

396 **Output:**
Two files are created and stored in

398 "~/Desktop/SeqCap/data/in_target_assemblies/":

400 1. "in_target.fasta": A fasta sequence file containing contigs that are stemmed from
the targeted loci.

402

2. "in_target.captured": A txt file containing percent captured for each target (tab-

404 delimited).

406 _____

408 *6-AssemblyEvaluation* (Optional): function "BASIC" evaluates the quality of in-
target assemblies by reporting basic stats: mean, median, total length, gc%, N50 etc.

410 It also generates a distribution of contigs by binned lengths.

412 In reality, "BASIC" can evaluates quality of any assemblies.

414 This script also assesses the quality of transcriptome/targeted capture assemblies
form several other aspects. For example:

416    "COVERAGE" calculates error rate, average quality score of the aligned bases and its variance/std, and average base coverage. Users need to generate alignment first.
418    "FIX" fixes assembly errors. Users need to generate alignment first.

420    For more details please execute "*6-AssemblyEvaluation*" in a terminal window.

422    *6-AssemblyEvaluation* BASIC
       **Input:**
424    "~/Desktop/SeqCap/data/in_target_assemblies/in_target.fasta" produced by *5-FindTargets*.
426

       **Commands:**
428    *ke@NGS:~/Desktop/SeqCap/data $ 6-AssemblyEvaluation BASIC -a ~/Desktop/SeqCap/data/in_target_assemblies*
430

       **Output:**
432    Two files are created in "~/Desktop/SeqCap/data/in_target_assemblies/":

434    1. "in_target.hist": distribution of contigs by binned lengths
       #Display first few lines of the file:
436    *ke@NGS:~/Desktop/SeqCap/data/in_target_assemblies$ head in_target.hist*
       *200:299       1026*
438    *300:399       242*
       *400:499       73*
440    *500:599       17*
       *600:699       7*
442    *700:799       0*
       *800:899       0*
444    *900:999       1*
       *1000:1099     0*
446    *1100:1199     1*

448    2. "basic_evaluation.out": results of assembly evaluation
       #Display first few lines of the file:
450    *ke@NGS:~/Desktop/SeqCap/data/in_target_assemblies$ head basic_evaluation.out*

452

       Note: users might want to compare the metric between in_target assemblies to
454    original targeted loci and see how much flanking are captured and assembled. Under most circumstances, the mean, median, N50 etc should be much higher in
456    in_target assemblies than in original targeted loci. The example dataset that is used for the purpose of demonstration, however, should
458    not show this pattern indicated above since it is assembled from a tiny fraction of data.
460    _____

462    *7-Alignment*: aligning cleaned reads against the pseudo-reference using Novoalign

464    Novoalign "is an aligner for single-ended and paired-end reads from the Illumina. Novoalign finds global optimum alignments using full Needleman-Wunsch
466    algorithm with affine gap penalties."

468    "Question: How does Novoalign compare to programs like BWA, Bowtie, ELAND and BFAST?
470    Answer:
Novoalign was designed to be an accurate short read aligner that combines fast K-
472    mer index searching with dynamic programming. In terms of speed Novoalign will be slower than Burrows-Wheeler transform aligners e.g. BWA, Bowtie and in some
474    cases faster than BFAST.  In terms of accuracy Novoalign is in most cases more sensitive than these tools because it uses full dynamic programming to find the best
476    alignment of a short read to a genome sequence."

478    According to Heng Li, author of SAMTools & MAQ, Novoalign "is the most accurate aligner to date".
480

Dependencies:
482    Novoalign: http://www.novocraft.com/main/downloadpage.php
SAMTools: http://sourceforge.net/projects/samtools/files/samtools/
484


486    **Input:**
1. A pseudo-reference genome, "~/Desktop/SeqCap/data/
488    in_target_assemblies/in_target.fasta", generated by *5-FindTargets*:
#make a new directory called "reference" under "~/Desktop/SeqCap/data/":
490    *ke@NGS:~/Desktop/SeqCap/data$ mkdir reference*

492    #Copy "in_target.fasta" to "~/Desktop/SeqCap/data/reference/":
*ke@NGS:~/Desktop/SeqCap/data$ cp in_target_assemblies/in_target.fasta reference/*
494

2. Cleaned reads generated by *2-ScrubReads:*
496    Cleaned reads are saved in "~/Desktop/SeqCap/data/cleaned_data/".
#Take a look at these reads:
498    *ke@NGS:~/Desktop/SeqCap/data/cleaned_data$ ls *.txt*
*CGRL_index1_1_final.txt*
500    *CGRL_index1_2_final.txt*
*CGRL_index1_u_final.txt*
502    *CGRL_index15_1_final.txt*
*CGRL_index15_2_final.txt*
504    *CGRL_index15_u_final.txt*
*CGRL_index40_1_final.txt*
506    *CGRL_index40_2_final.txt*
*CGRL_index40_u_final.txt*

11

508

**Commands:**

510

#Make a new folder called "alignment" under "~/Desktop/SeqCap/data/":
512 *ke@NGS:~/Desktop/SeqCap/data$ mkdir alignment*

514 #Run *7-Alignment:*
*ke@NGS:~/Desktop/SeqCap/data$ 7-Alignment -f*
516 *~/Desktop/SeqCap/data/reference/in_target.fasta -r*
*~/Desktop/SeqCap/data/cleaned_data/ -o ~/Desktop/SeqCap/data/alignment/ -i*
518 *200 -v 20 -t 90*

520 Note: do not set t for alignment of very divergent genomes.

522 **Output:**
BAMS and indexed bam files.
524 #Take a look at these files:
*ke@NGS:~/Desktop/SeqCap/data/alignment$ ls*
526 *CGRL_index1_sorted.bam*
*CGRL_index1_sorted.bam.bai*
528 *CGRL_index15_sorted.bam*
*CGRL_index15_sorted.bam.bai*
530 *CGRL_index40_sorted.bam*
*CGRL_index40_sorted.bam.bai*
532 _____

534 *8-ExonCaptureEvaluation* (Optional): Function "Evaluation" provides evaluation
for capture efficiency: %reads mapped, %target captured, average sequence depth,
536 etc.

538 Note: %reads mapped (specificity), %target captured (sensitivity), and average
sequence depth are typically reported in papers.
540
Dependencies:
542 SAMTools: http://sourceforge.net/projects/samtools/files/samtools/
BEDTools: http://bedtools.readthedocs.org/en/latest/content/installation.html
544

*8-ExonCaptureEvaluation* Evaluation
546 **Input:**
1. A pseudo-reference "in_target.fasta" generated by *5-FindTargets*:
548 You can find this file in "~/Desktop/SeqCap/data/reference/".

550 2. Cleaned reads generated by *2-ScrubReads*:
These reads are located in "~/Desktop/SeqCap/data/cleaned_data/":
552 CGRL_index1_1_final.txt
CGRL_index1_2_final.txt

| | |
|---|---|
| 554 | CGRL_index1_u_final.txt |
| | CGRL_index15_1_final.txt |
| 556 | CGRL_index15_2_final.txt |
| | CGRL_index15_u_final.txt |
| 558 | CGRL_index40_1_final.txt |
| | CGRL_index40_2_final.txt |
| 560 | CGRL_index40_u_final.txt |

562   3. Raw reads generated by *1-PreCleanup*:
These data are located in "~/Desktop/SeqCap/data/rawdata/raw/pre-clean/":
564   CGRL_index1_R1.fq
CGRL_index1_R2.fq
566   CGRL_index15_R1.fq
CGRL_index15_R2.fq
568   CGRL_index40_R1.fq
CGRL_index40_R2.fq
570

4. All bam (alignment) files generated by *7-Alignment*:
572   The bams (sorted and indexed) are located in
"~/Desktop/SeqCap/data/alignment/":
574   CGRL_index1_sorted.bam
CGRL_index1_sorted.bam.bai
576   CGRL_index15_sorted.bam
CGRL_index15_sorted.bam.bai
578   CGRL_index40_sorted.bam
CGRL_index40_sorted.bam.bai
580

5. A .bed file generated by *9-preFiltering* (optional)
582   A BED file (.bed) is a tab-delimited text file that defines a feature track of each locus.
In this case, this file defines targeted region in each assembled contig.
584

For example if the length of contig125 is 1000bp, but the targeted region starts from
586   position 120 and ends by 350, then the correct expression is:

588   Contig125 119 350  (note: in bed the start position is one less than it's actual value)

590   For more details of BED format please go to:
http://www.broadinstitute.org/igv/BED
592

**Commands:**
594   #Make a new folder called "ExonCapEval" under "~/Desktop/SeqCap/data/":
*ke@NGS:~/Desktop/SeqCap/data$ mkdir ExonCapEval*
596

#Run *8-ExonCaptureEvaluation*:
598   *ke@NGS:~/Desktop/SeqCap/data$ 8-ExonCaptureEvaluation Evaluation  -genome
~/Desktop/SeqCap/data/reference/in_target.fasta -cleanDir*

600     *~/Desktop/SeqCap/data/cleaned_data/ -rawDir*
        *~/Desktop/SeqCap/data/rawdata/raw/pre-clean/ -bamDir*
602     *~/Desktop/SeqCap/data/alignment/ -InstrID HS -resDir*
        *~/Desktop/SeqCap/data/ExonCapEval/ -readlen 100*
604
        Note: If you just evaluate how targeted regions worked, you should provide a bed
606     file (generated by 9-preFiltering) while running 8-ExonCaptureEvaluation
        Evaluation.
608
        **Output:**
610     "data_metrics.txt" under "~/Desktop/SeqCap/data/ExonCapEval/"
        You can use "less" to check the results reported in this file.
612

614     _____

        *\*9-preFiltering\**: "*9-preFiltering* bed" generates a bed for exonic region(s) from each
616     contig in in-target assemblies (aka. the reference) and a bed for all assembled
        contigs (start position is 0 and the end position is the length of the contig); "*9-*
618     *preFiltering percentile*" produces a list of contigs that fall outside the desired
        coverage percentiles; "*9-preFiltering percentile*" also produces base coverage values
620     at different level of percentile.

622     Dependencies:
        Tie-Array-Packed-0.13: http://search.cpan.org/~salva/Tie-Array-Packed-
624     0.13/lib/Tie/Array/Packed.pm

626     *9-preFiltering bed:*
        **Input:**
628     1.Targeted loci:
        "~/Desktop/SeqCap/denovoTargetCapture/original_target/targeted_loci.fasta";
630
        2. "~/Desktop/SeqCap/data/reference/in_target.fasta" generated by *5-FindTargets*.
632
        **Commands:**
634     #Make a new folder called "bed_files" under "~/Desktop/SeqCap/data/":
        *ke@NGS:~/Desktop/SeqCap/data$ mkdir bed_files*
636
        #cd to this folder:
638     *ke@NGS:~/Desktop/SeqCap/data$ cd bed_files*

640     #Run *9-preFiltering* bed:
        *ke@NGS:~/Desktop/SeqCap/data/bed_files$ 9-preFiltering bed*
642     *~/Desktop/SeqCap/denovoTargetCapture/original_target/targeted_loci.fasta*
        *~/Desktop/SeqCap/data/reference/in_target.fasta*
644

14

**Output**:

646 Two file under "~/Desktop/SeqCap/data/bed_files/":

1. "final.bed" is used as input for *9-preFiltering* percentile and *8-*

648 *ExonCaptureEvaluation* Evaluation

2. "All_contig.bed" is used as input for *10-SNPcleaner*.

650


652 *9-preFiltering* percentile:

**Input:**

654 1.  Make a new folder called "pre-filtering" under "~/Desktop/SeqCap/data/" and cd
to this folder:

656 *ke@NGS:~/Desktop/SeqCap/data$ mkdir pre-filtering*

*ke@NGS:~/Desktop/SeqCap/data$ cd pre-filtering*

658

2. In "~/Desktop/SeqCap/data/pre-filtering/", generate a merged, sorted bam for

660 all samples:

*ke@NGS:~/Desktop/SeqCap/data/pre-filtering$ samtools merge merge.bam*

662 *~/Desktop/SeqCap/data/alignment/*.bam*

*ke@NGS:~/Desktop/SeqCap/data/pre-filtering$ samtools sort merge.bam*

664 *merge_sorted*


666 "~/Desktop/SeqCap/data/pre-filtering/merge_sorted.bam" is the input bam.


668 3. A bed file:

"~/Desktop/SeqCap/data/bed_files/final.bed" is generated by *9-preFiltering* bed

670

**Commands:**

672 # Run *9-preFiltering* percentile under "~/Desktop/SeqCap/data/pre-filtering/":

*ke@NGS:~/Desktop/SeqCap/data/pre-filtering$ 9-preFiltering percentile -b*

674 *~/Desktop/SeqCap/data/pre-filtering/merge_sorted.bam -o CGRL -B*

*~/Desktop/SeqCap/data/bed_files/final.bed*

676

**Output**:

678 In the folder "~/Desktop/SeqCap/data/pre-filtering/" there are a couple of files
created:

680 1. "CGRL_gene_outside_percentile.txt" shows a list of contigs having coverage <X%
or >Y% percentiles of the data. X and Y are defined by users. This file will be used in

682 *10-SNPcleaner*.

2. "CGRL_site_depth_percentile.txt" shows base coverage at different level of

684 percentiles.  The information in this file will be used by *10-SNPcleaner*.

3. "CGRL_gene_depth_percentile.txt" shows average base coverage at different level

686 of percentiles.

4. "CGRL_gene_depth.txt" shows average coverage of each contig. If you want to

688 know more about empirical coverage distribution of your data then you take the
coverage value from this file and use R to plot it.

690    5. "CGRL_site_depth.txt" shows per-base coverage of the data. You can plot it to get a
       sense of empirical distribution of base coverage.
692    6. "CGRL_gene_outside_sd_filter.txt": shows a list of contigs all outside N standard
       deviation of the mean. Users set N when running the command.
694

696    Note: users might want to perform filtering based on other criteria such as 3
       standard deviations of the mean. However, this method usually requires a normal
698    distribution of the data. In reality per-base depth of exon capture data rarely follows
       a normal distribution.
700    _____

702    *10-SNPcleaner*: Raw variant filtering and generates a "keep" file for the following
       SNP/genotype calling by ANGSD. This script is mainly for filtering data at contig and
704    site levels. Users need to perform individual-level filtering before running this
       script. See below for more details.
706
       Before we call SNPs /genotypes and estimate allele frequencies using ANGSD, we
708    usually employ three levels of filtering on the data sets in a hierarchical order:
       individual level, contig level and site level. The filters in each step of the hierarchy
710    are applied only to the subset of data that pass the quality control thresholds at all
       previous levels. The first filters applied are the individual-level filters to remove
712    entire individuals deviating excessively from the average across-individual coverage
       and error rate. Contig-level filters, followed by site-level filters, are then applied to
714    remove entire contigs and sites, respectively, that appeared to be quality outliers.
       All individual specimens, contigs and sites should be filtered on multiple aspects of
716    quality (e.g. potential cross-sample DNA contamination, sequencing errors,
       paralogy).
718
       1. Filtering at individual level
720    a. Remove individuals having extremely low or high coverage. Individual coverage
       can be estimated using *8-ExonCaptureEvaluation Evaluation.* The file you want to
722    examine is "~/Desktop/SeqCap/data/ExonCapEval/data_metrics.txt"

724    *ke@NGS:~/Desktop/SeqCap/data$ less –S ExonCapEval/data_metrics.txt*

726    b. Remove individuals with excessively high sequencing error rates measured as the
       percentage of mismatched bases out of the total number of aligned bases in the
728    mitochondrial genome. Empirical error can be estimated using *6-
       AssemblyEvaluation COVERAGE*
730
       To run *6-AssemblyEvaluation COVERAGE* you need first to generate pileup files for
732    mitochondrial locus for each sample.

734    *ke@NGS:~/Desktop/SeqCap/data$ 6-AssemblyEvaluation COVERAGE*

736    *Usage 6-AssemblyEvaluation COVERAGE [options]*

738    *Options:*
       *-p   DIR      folder containing all pileup*
740    *                files generated by "samtools*
       *                mpileup -f ref.fa sample1.bam*
742    *                 > sample1.pileup"*
       *-c   INT      coverage cutoff [5]*
744    *-q   INT       base quality cutoff [13]*

746

## 2. Filtering at contig level

748    a. Remove contigs that show extremely low or high coverage based on the empirical
       coverage distribution across all contigs. *9-preFiltering* *percentile* can be used to
750    generate a list of contigs that show extreme coverage based on percentile values (for
       example: 1% and 99%; 5% and 95% etc.). This list can then be used as one of input
752    files in *10-SNPcleaner* for the purpose of filtering.

754    b. Remove contigs with at least one SNP having allele frequencies highly deviating
       from Hardy–Weinberg equilibrium expectations. Done by *10-SNPcleaner*. Note this
756    is a very stringent filter even for exon capture dataset and not suitable at all for
       genomic dataset. To use this filter you need to provide
758    "~/Desktop/SeqCap/data/bed_files/All_contig.bed" generated by *9-preFiltering* *bed.*

## 3. Filtering at site level
760    a. Remove sites with excessively low or high coverage based on the empirical
       coverage distribution. To determine high (e.g. 99% or 95%) and low (e.g. 1% or 5%)
762    percentiles of base coverage you need run *9-preFiltering* *percentile* to get
       "CGRL_site_depth_percentile.txt".
764

766    b. Remove sites having allele frequencies highly deviating from Hardy–Weinberg
       equilibrium expectations (exact test). Done by *10-SNPcleaner*. This filter can be
768    combined with the contig HWE filter (2.b).

770    c. Remove sites with biases associated with reference and alternative allele Phred
       quality, mapping quality and distance of alleles from the ends of reads. Also remove
772    sites that show a bias towards SNPs coming from the forward or reverse strand.
       These will be done by *10-SNPcleaner*.
774
       -> Strand Bias: Tests if variant bases tend to come from one strand.
776    -> End Distance Bias: Tests if variant bases tend to occur at a fixed distance from the
       end of reads, which is usually an indication of misalignment.
778    -> Base Quality Bias: Tests if variant bases tend to occur with a Phred-scale quality
       bias.
780    -> Mapping Quality Bias: Tests if variant bases tend to occur with a mapping quality
       bias.

782

d. Remove sites for which there are not at least M of the individuals sequenced at N
784 coverage each. This makes sure that the remaining data matrix does not contain too
much missing data. This will be done by *10-SNPcleaner*.
786

e. Remove sites with a root mean square (RMS) mapping quality for SNPs across all
788 samples below a certain threshold. It is a measure of the variance of quality scores.
This will be done by *10-SNPcleaner*.
790

f. (optional) For historic samples, characterize the pattern of base mis-incorporation
792 first. Sometimes it is necessary to remove C to T and G to A SNPs from the dataset.
This can be done by *10-SNPcleaner*.
794

Before running *10-SNPcleaner* make sure that individual-level filtering is finished.
796

**Input**:
798 1. "~/Desktop/SeqCap/data/pre-filtering/CGRL_gene_outside_percentile.txt" by *9-preFiltering* percentile.
800

2. "~/Desktop/SeqCap/data/bed_files/All_contig.bed" generated by *9-preFiltering*
802 *bed*.

804 3. "~/Desktop/SeqCap/data/pre-filtering/CGRL_site_depth_percentile.txt" by *9-preFiltering* percentile* is ready and will be used to guide the site-level coverage
806 filtering.

808 4.  Create a new folder "SNPcleaning" under "~/Desktop/SeqCap/data/" and inside
this folder generate a raw vcf that contains all sites from all individual samples that
810 pass individual-level filters:
*ke@NGS:~/Desktop/SeqCap/data$ mkdir SNPcleaning*
812 *ke@NGS:~/Desktop/SeqCap/data$ cd SNPcleaning/*
*ke@NGS:~/Desktop/SeqCap/data/SNPcleaning$ samtools mpileup -B -D -I -S  -uf*
814 *~/Desktop/SeqCap/data/reference/in_target.fasta*
*~/Desktop/SeqCap/data/alignment/*sorted.bam | bcftools view -cg - > raw.vcf*
816

-D       output per-sample DP in BCF
818 -B        disable BAQ computation
-I        do not perform indel calling
820 -S        output per-sample strand bias P-value in BCF

822 **Commands:**
#Run *10-SNPcleaner* under "~/Desktop/SeqCap/data/SNPcleaning/":
824 *ke@NGS:~/Desktop/SeqCap/data/SNPcleaning$ 10-SNPcleaner -d 2 -D 7 -k 2 -u 1 -a 0*
*-B CGRL.bed -p CGRL_filtered -r ~/Desktop/SeqCap/data/pre-*
826 *filtering/CGRL_gene_outside_percentile.txt -X*
*~/Desktop/SeqCap/data/bed_files/All_contig.bed -g -v raw.vcf> out.vcf*

828

Note: for "-D 7", 7 is the 99% percentile of the base coverage. We get this number
830   from "~/Desktop/SeqCap/data/pre-filtering/CGRL_site_depth_percentile.txt".

832   **Output:**
In "~/Desktop/SeqCap/data/SNPcleaning/", several files are created:
834

1. "CGRL.bed" contains sites  (potentially variable and non-variable) passing all
836   filters.

838   #To generate a keep file for ANGSD:
*ke@NGS:~/Desktop/SeqCap/data/SNPcleaning$ cut -f1,2 CGRL.bed > CGRL.keep*
840

2. "CGRL_filtered" (dumped with option -p) contains all sites that failed to pass
842   certain filters. Characters in front of filtered sites indicate filters that the site failed
to pass.
844

#To view this file:
846   *ke@NGS:~/Desktop/SeqCap/data/SNPcleaning$ bunzip2 -c CGRL_filtered  | less -S*

848   3. "out.vcf" is the resulting vcf that contains sites (both variable and non-variable)
passed all filters
850

Questions:
852   1. Check how many sites are present before and after filtering?
2. Check why some sites are filtered out by examining "CGRL_filtered".
854