

Pipeline for *de novo* Targeted Capture

PART I population genomics

March 5, 2015

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/CGRL-QB3-UCBerkeley/denovoTargetCapturePopGen>

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](#)

[11-PopGenTools](#)

Use “`chmod +x script`” to make each of these perl scripts executable.

52

54 **1-PreCleanup**: Reformats raw sequencing reads from Illumina HiSeq or MiSeq for
56 [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:
FastQC: <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

60

Input:

62 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
64 saved under “/home/ke/Desktop/SeqCap/data/rawdata/library/”. Compressed
fastq sequence files are saved in each of these folders.

66

Fastq files use the following naming scheme:

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

70

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

72 sample name: CGRL_index15
barcode sequence: CGACCTG
74 lane (0-padded to 3 digits): 006
read number: 1
76 set number (0-padded to 3 digits): 001

78

#Make a new folder called “raw” under “~/Desktop/SeqCap/data/rawdata/”:

80 `ke@NGS:~/Desktop/SeqCap/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/SeqCap/data/rawdata$ cp library/CGRL_index*/*.gz raw/`

86

#Check data files in “raw”:

88 `ke@NGS:~/Desktop/SeqCap/data/rawdata$ ls raw/*`

90

92 **Commands:**
#cd to the working directory:

```

94  ke@NGS:~/Desktop/SeqCap/data/rawdata$ cd ..

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

100
Output:
102  Three new folders will be created under "~/Desktop/SeqCap/data/rawdata/raw/":
    "pre-clean"
104  "combined"
    "pre-clean/evaluation"
106
    - Folder "pre-clean" contains reformatted raw fastq reads.
108  CGRL_index1_R1.fq
    CGRL_index1_R2.fq
110  CGRL_index15_R1.fq
    CGRL_index15_R2.fq
112  CGRL_index40_R1.fq
    CGRL_index40_R2.fq
114
    - Folder "combined" contains merged, compressed, fastq data files (not used by the
116  following pipeline).
    CGRL_index1_TCGCAGG_L006_R1.fastq.gz
118  CGRL_index1_TCGCAGG_L006_R2.fastq.gz
    CGRL_index15_CGACCTG_L006_R1.fastq.gz
120  CGRL_index15_CGACCTG_L006_R2.fastq.gz
    CGRL_index40_TTCGCAA_L006_R1.fastq.gz
122  CGRL_index40_TTCGCAA_L006_R2.fastq.gz

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

132  Questions:
134  1. Check the sequence identifiers and the number of reads in fastq files before and
    after running 1-PreCleanup and compare the results.
    2. Check the fastQC evaluation results for the raw data
136

```

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

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

152 **Input:**

1. Reformatted fastq files created by [1-PreCleanup](#):
154 #Check the raw data files:

```
ke@NGS:~/Desktop/SeqCap/data/rawdata/raw/pre-clean$ ls *.fq
```

156 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 AATGATACGGCGACCACCGAGATCTACACcctgcgaACACTCTTTCCCTACACGACGCTCTTCCGATCT  
>P5_index2  
176 AATGATACGGCGACCACCGAGATCTACACtgcagagACACTCTTTCCCTACACGACGCTCTTCCGATCT  
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

```
library      P7    P5
```

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 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/SeqCap/data/”:

```
ke@NGS:~/Desktop/SeqCap/data$ mkdir cleaned_data
```

204

#Run [2-ScrubReads](#):

206

```
ke@NGS:~/Desktop/SeqCap/data$ 2-ScrubReads -f
```

```
~/Desktop/SeqCap/data/rawdata/raw/pre-clean/ -o
```

208

```
~/Desktop/SeqCap/data/cleaned_data/ -a
```

```
~/Desktop/SeqCap/denovoTargetCapture/associated_files/Adapters.fasta -b
```

210

```
~/Desktop/SeqCap/denovoTargetCapture/associated_files/libInfo.txt -t
```

```
/home/ke/Desktop/SeqCap/programs/Trimmomatic-0.32/trimmomatic-0.32.jar -c
```

212

```
~/Desktop/SeqCap/denovoTargetCapture/associated_files/ecoli/e_coli_K12.fasta -e  
200 -z
```

214

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

218

Output:

1. In “~/Desktop/SeqCap/data/cleaned_data/”, six .txt files per library are produced:

220

For example for library CGRL_index1, the six files are:

222

CGRL_index1_1_final.txt (left reads)

CGRL_index1_2_final.txt (right reads)

224

CGRL_index1_u_final.txt (merged or unpaired reads)

CGRL_index1.contam.out (headers of reads aligned to bacteria)

226

CGRL_index1.duplicates.out (headers of duplicated reads)

CGRL_index1.lowComplexity.out (headers of low complexity reads)

228

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

230

232 Questions:

1. Check the %reads that are exact duplicates, %reads that are likely derived from microbial genome and %reads that contain low complexity.

234

2. Check the fastQC evaluation results of cleaned reads and then compare them to those of raw reads. Is the quality improved?

236

238

240 **3-GenerateAssemblies**: Assemble sequence capture data using ABySS.

242 We use a multi-kmer approach to assemble our data. If there is even coverage and
 244 even polymorphism levels across the assembled genome, there should (in theory)
 246 be one k-mer that best assembles the data. In reality, coverage and polymorphism
 248 vary across captured loci, and using multiple k-mers is a way to bet hedge and get
 250 good assemblies for all loci. In assembling your data, it is important to consider
 252 which samples to use in your assembly. Ideally, you could assemble across multiple
 254 individuals to increase your read depth, and thus, assembly contiguity and
 256 continuity. However, for many projects, more individuals can also mean increased
 258 polymorphism. While we have found the assemblers are more robust to
 260 polymorphism than the program writers themselves often suggest, increased
 262 polymorphism does lead to shorter contigs and increased misassemblies. With these
 264 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
 268 much polymorphism.

Dependencies:

258 ABySS (compiled with OpenMPI and Google sparsehash):
<http://www.bcgsc.ca/platform/bioinfo/software/abyss>
 260 ABySS uses OpenMP for parallelization

262 **Input:**
 Concatenated cleaned reads from libraries that you would like to assemble together.
 264 The libraries to be assembled together have to be genetically similar: ideally,
 samples from the same population. In this example we want to assemble
 266 CGRL_index1, CGRL_index15 and CGRL_index40 together.

268 #Make a new folder called “raw_assembly” under “~/Desktop/SeqCap/data/”:
 ke@NGS:~/Desktop/SeqCap/data\$ *mkdir raw_assembly*

270 #Concatenate cleaned reads and save them in “raw_assembly”:

272 ke@NGS:~/Desktop/SeqCap/data\$ *cat cleaned_data/CGRL_index*_1_final.txt >*
raw_assembly/combined_1_final.txt
 274 ke@NGS:~/Desktop/SeqCap/data\$ *cat cleaned_data/CGRL_index*_2_final.txt >*
raw_assembly/combined_2_final.txt
 276 ke@NGS:~/Desktop/SeqCap/data\$ *cat cleaned_data/CGRL_index*_u_final.txt >*
raw_assembly/combined_u_final.txt

278

280 #Inside “raw_assemblies” make a new folder “results”:
 ke@NGS:~/Desktop/SeqCap/data\$ *mkdir raw_assembly/results*

282

284

Commands:

286 #Run ABySS on two processors using kmer sizes of 21, 31, 41, 51, 61, and 71 (do not
execute the command in the workshop).

```
288 ke@NGS:~/Desktop/SeqCap/data$ 3-GenerateAssemblies abyss -reads  
~/Desktop/SeqCap/data/raw_assembly/ -mpi /usr/bin/mpirun -out  
290 ~/Desktop/SeqCap/data/raw_assembly/results/ -kmer 21 31 41 51 61 71 -np 2
```

292 **Note: Your laptop will not be able to handle memory intensive ABySS assemblies.
Please do not run it at the class. You can copy the result files saved in the back-up
294 folder to “~/Desktop/SeqCap/data/raw_assembly/results/”.

```
ke@NGS:~/Desktop/SeqCap/data$ scp -r  
296 ~/Desktop/SeqCap/denovoTargetCapture/associated_files/combined/  
/home/ke/Desktop/SeqCap/data/raw_assembly/results
```

298

300

Output:

302 There are a lot of intermediate files created in
“~/Desktop/SeqCap/data/raw_assembly/results/combined/”.

304

#To show the assemblies that we need for the next step:

```
306 ke@NGS:~/Desktop/SeqCap/data$ cd raw_assembly/results/combined/  
ke@NGS:~/Desktop/SeqCap/data/raw_assembly/results/combined$ ls *-contigs.fa
```

308 combined_k21_cov_default-contigs.fa

combined_k31_cov_default-contigs.fa

310 combined_k41_cov_default-contigs.fa

combined_k51_cov_default-contigs.fa

312 combined_k61_cov_default-contigs.fa

combined_k71_cov_default-contigs.fa

314

#Combine all the raw assemblies and write the result to a new file called
316 “all_assemblies.fasta”:

```
ke@NGS:~/Desktop/SeqCap/data/raw_assembly/results/combined$ cat  
318 combined_*_cov_default-contigs.fa > all_assemblies.fasta
```

320 #Make a new folder called “merge_assemblies” under “~/Desktop/SeqCap/data/”:

```
ke@NGS:~/Desktop/SeqCap/data $ mkdir merge_assemblies
```

322

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

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

326

328

330

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

336 Dependencies:

CAP3: <http://seq.cs.iastate.edu/cap3.html>

338 blat: <http://users.soe.ucsc.edu/~kent/src/>

cd-hit-est: <https://code.google.com/p/cdhit/downloads/list>

340

Input:

342 Concatenated raw assemblies

344 “~/Desktop/SeqCap/data/merge_assemblies/all_assemblies.fasta” produced by [3-GenerateAssemblies](#)

346 **Commands:**

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

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

352

Output:

354 Several files are created in “~/Desktop/SeqCap/data/merge_assemblies/”. The
356 data file that we need for the next step is “all_assemblies.fasta.final”.

356

358

360 **5-FindTargets**: identify contigs that are stemmed from the targeted loci and use
 362 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
 368 more contigs than loci you targeted. Some of these might be junk; some might be
 370 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
 implement a BLAST approach, which identifies which contig has the best-hit match
 to one's targeted loci.

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

376 **Input:**
 1. "~/Desktop/SeqCap/data/merge_assemblies/all_assemblies_final.fasta"
 378 produced by [4-FinalAssembly](#).

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

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

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

388 #Copy the merged assemblies to "in_target_assemblies":

```
ke@NGS:~/Desktop/SeqCap/data$ cp merge_assemblies/all_assemblies.fasta.final  
in_target_assemblies/
```

392

394 #Run [5-FindingTargets](#) in "~/Desktop/SeqCap/data/in_target_assemblies/" (do not
 run this script in the class):

```
ke@NGS:~/Desktop/SeqCap/data$ 5-FindingTargets -t  
~/Desktop/SeqCap/denovoTargetCapture/original_target/targeted_loci.fasta -a  
398 ~/Desktop/SeqCap/data/in_target_assemblies/
```

400 #Copy the results from the backup folder to
 "~/Desktop/SeqCap/data/in_target_assemblies/"

```
402 ke@NGS:~/Desktop/SeqCap/data$ scp -r  
    ~/Desktop/SeqCap/denovoTargetCapture/associated_files/In_target/  
404 in_target_assemblies/
```

406

Output:

408 Output files are created and stored in
 “~/Desktop/SeqCap/data/in_target_assemblies/In_target”:

410

412 1. “all_assemblies_targetedRegionAndFlanking.fasta”: A fasta sequence file
 containing contigs that are stemmed from the targeted loci (exons + flanking).

414 2. “all_assemblies_intarget.sensitivity”: A txt file containing percent captured for
 each target (tab-delimited).

416

418 3. “all_assemblies_targetedRegionAndFlanking.bed”: A bed file contains coordinates
 of the targeted regions. This bed file will be used in [8-ExonCaptureEvaluation](#)
 Evaluation and [9-preFiltering](#).

420

422 4. “all_assemblies_contig.bed”: A bed file contains coordinates of the
 targeted+flanking regions. This bed file will be used in [10-SNPcleaner](#)

424

426 *6-AssemblyEvaluation* (Optional): function “BASIC” evaluates the quality of in-
428 target assemblies by reporting basic stats: mean, median, total length, gc%, N50 etc.
It also generates a distribution of contigs by binned lengths.

430 In reality, “BASIC” can evaluates quality of any assemblies.

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

434 “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.

436 “FIX” fixes assembly errors. Users need to generate alignment first.

438 For more details please execute “[6-AssemblyEvaluation](#)” in a terminal window.

440 [6-AssemblyEvaluation](#) BASIC

Input:

442 “~/Desktop/SeqCap/data/in_target_assemblies/in_target.fasta” produced by [5-
FindTargets](#).

444

Commands:

446 ke@NGS:~/Desktop/SeqCap/data\$ **6-AssemblyEvaluation BASIC -a
~/Desktop/SeqCap/data/in_target_assemblies/In_target/**

448

450

Output:

452 Two files are created in “~/Desktop/SeqCap/data/in_target_assemblies/In_target”:

454 1. “in_target.hist”: distribution of contigs by binned lengths

#Display first few lines of the file:

456 ke@NGS:~/Desktop/SeqCap/data/in_target_assemblies/In_target\$ **head in_target.hist**

	200:299	1026
458	300:399	242
	400:499	73
460	500:599	17
	600:699	7
462	700:799	0
	800:899	0
464	900:999	1
	1000:1099	0
466	1100:1199	1

468 2. "basic_evaluation.out": results of assembly evaluation
#Display the content of the file:
470 *ke@NGS:~/Desktop/SeqCap/data/in_target_assemblies/In_target\$ less -S*
basic_evaluation.out
472
474 Note: users might want to compare the metric between in_target assemblies to
original targeted loci and see how much flanking are captured and assembled.
476 Under most circumstances, the mean, median, N50 etc should be much higher in
in_target assemblies than in original targeted loci.
478 The example dataset that is used for the purpose of demonstration, however, should
not show this pattern indicated above since it is assembled from a tiny fraction of
480 data.

482

484 **7-Alignment**: aligning cleaned reads against the pseudo-reference using Novoalign

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

490 “Question: How does Novoalign compare to programs like BWA, Bowtie, ELAND and
BFAST?

492 Answer:

Novoalign was designed to be an accurate short read aligner that combines fast K-
494 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
496 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
498 alignment of a short read to a genome sequence.”

500 According to Heng Li, author of SAMTools & MAQ, Novoalign “is the most accurate
aligner to date”.

502

Dependencies:

504 Novoalign: <http://www.novocraft.com/main/downloadpage.php>

SAMTools: <http://sourceforge.net/projects/samtools/files/samtools/>

506

508 **Input:**

1. A pseudo-reference genome, “~/Desktop/SeqCap/data/
510 in_target_assemblies/in_target.fasta”, generated by [5-FindTargets](#):
#make a new directory called “reference” under “~/Desktop/SeqCap/data/”:

512 *ke@NGS:~/Desktop/SeqCap/data\$ mkdir reference*

514

#Copy “*all_assemblies_targetedRegionAndFlanking.fasta*” to
516 “~/Desktop/SeqCap/data/reference/”:

ke@NGS:~/Desktop/SeqCap/data\$ cp
518 *in_target_assemblies/In_target/all_assemblies_targetedRegionAndFlanking.fasta*
reference/

520

522 2. Cleaned reads generated by [2-ScrubReads](#):
Cleaned reads are saved in “~/Desktop/SeqCap/data/cleaned_data/”.

524 #list all the reads:


```
ke@NGS:~/Desktop/SeqCap/data/cleaned_data$ ls *.txt
```

```
CGRL_index1_1_final.txt  
CGRL_index1_2_final.txt  
CGRL_index1_u_final.txt  
CGRL_index15_1_final.txt  
CGRL_index15_2_final.txt  
CGRL_index15_u_final.txt  
CGRL_index40_1_final.txt  
CGRL_index40_2_final.txt  
CGRL_index40_u_final.txt
```

Commands:

```
#Make a new folder called "alignment" under "~/Desktop/SeqCap/data/":
```

```
ke@NGS:~/Desktop/SeqCap/data$ mkdir alignment
```

```
#Run 7-Alignment:
```

```
ke@NGS:~/Desktop/SeqCap/data$ 7-Alignment -f  
~/Desktop/SeqCap/data/reference/all_assemblies_targetedRegionAndFlanking.fasta  
-r ~/Desktop/SeqCap/data/cleaned_data/ -o ~/Desktop/SeqCap/data/alignment/ -i  
200 -v 20 -t 90
```

```
Note: do not set "-t" for alignment to very divergent ref genomes.
```

Output:

```
BAMS and indexed bam files.
```

```
#Take a look at these files:
```

```
ke@NGS:~/Desktop/SeqCap/data/alignment$ ls
```

```
CGRL_index1_sorted.bam  
CGRL_index1_sorted.bam.bai  
CGRL_index15_sorted.bam  
CGRL_index15_sorted.bam.bai  
CGRL_index40_sorted.bam  
CGRL_index40_sorted.bam.bai
```

566 *8-ExonCaptureEvaluation* (Optional): Function “Evaluation” provides evaluation
568 for capture efficiency: %reads mapped, %target captured, average sequence depth,
etc.

570 Note: %reads mapped (specificity), %target captured (sensitivity), and average
572 sequence depth are typically reported in papers.

Dependencies:

574 SAMTools: <http://sourceforge.net/projects/samtools/files/samtools/>

BEDTools: <http://bedtools.readthedocs.org/en/latest/content/installation.html>

576

[8-ExonCaptureEvaluation](#) Evaluation

578 **Input:**

1. A pseudo-reference “all_assemblies_targetedRegionAndFlanking.fasta” generated
580 by [5-FindTargets](#):

You can find this file in “~/Desktop/SeqCap/data/reference/”.

582

2. Cleaned reads generated by [2-ScrubReads](#):

584 These reads are located in “~/Desktop/SeqCap/data/cleaned_data/”:

CGRL_index1_1_final.txt

586 CGRL_index1_2_final.txt

CGRL_index1_u_final.txt

588 CGRL_index15_1_final.txt

CGRL_index15_2_final.txt

590 CGRL_index15_u_final.txt

CGRL_index40_1_final.txt

592 CGRL_index40_2_final.txt

CGRL_index40_u_final.txt

594

3. Raw reads generated by [1-PreCleanup](#):

596 These data are located in “~/Desktop/SeqCap/data/rawdata/raw/pre-clean/”:

CGRL_index1_R1.fq

598 CGRL_index1_R2.fq

CGRL_index15_R1.fq

600 CGRL_index15_R2.fq

CGRL_index40_R1.fq

602 CGRL_index40_R2.fq

604 4. All bam (alignment) files generated by [7-Alignment](#):

The bams (sorted and indexed) are located in

606 “~/Desktop/SeqCap/data/alignment/”:

CGRL_index1_sorted.bam

608 CGRL_index1_sorted.bam.bai

CGRL_index15_sorted.bam

610 CGRL_index15_sorted.bam.bai

CGRL_index40_sorted.bam
 612 CGRL_index40_sorted.bam.bai

614 5. A .bed file “all_assemblies_targetedRegionAndFlanking.bed” in the folder
 “~/Desktop/SeqCap/data/in_target_assemblies/In_target/”
 616 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.
 618
 For example if the length of contig125 is 1000bp, but the targeted region starts from
 620 position 120 and ends by 350, then the correct expression is:
 622 Contig125 119 350 (note: in bed the start position is one less than it's actual value)
 624 For more details of BED format please go to:
<http://www.broadinstitute.org/igv/BED>
 626

Commands:

628 #Make a new folder called “ExonCapEval” under “~/Desktop/SeqCap/data/”:
 ke@NGS:~/Desktop/SeqCap/data\$ *mkdir ExonCapEval*

630

632 #Run [8-ExonCaptureEvaluation](#) without a bed file:
 ke@NGS:~/Desktop/SeqCap/data\$ *8-ExonCaptureEvaluation Evaluation -genome*
 634 *~/Desktop/SeqCap/data/reference/in_target.fasta -cleanDir*
~/Desktop/SeqCap/data/cleaned_data/ -rawDir
 636 *~/Desktop/SeqCap/data/rawdata/raw/pre-clean/ -bamDir*
~/Desktop/SeqCap/data/alignment/ -InstrID HS -resDir
 638 *~/Desktop/SeqCap/data/ExonCapEval/ -readlen 100*

640 #Run [8-ExonCaptureEvaluation](#) with a bed file:
 ke@NGS:~/Desktop/SeqCap/data\$ *8-ExonCaptureEvaluation Evaluation -genome*
 642 *~/Desktop/SeqCap/data/reference/in_target.fasta -cleanDir*
~/Desktop/SeqCap/data/cleaned_data/ -rawDir
 644 *~/Desktop/SeqCap/data/rawdata/raw/pre-clean/ -bamDir*
~/Desktop/SeqCap/data/alignment/ -InstrID HS -resDir
 646 *~/Desktop/SeqCap/data/ExonCapEval/ -readlen 100 -bedFile*
in_target_assemblies/In_target/all_assemblies_targetedRegionAndFlanking.bed

648

Output:
 650 “data_metrics.txt” under “~/Desktop/SeqCap/data/ExonCapEval/”
 You can use “less” to check the results reported in this file.
 652

654 *9-preFiltering*: “[9-preFiltering percentile](#)” produces a list of contigs that fall outside
656 the desired coverage percentiles; “[9-preFiltering percentile](#)” also produces base
658 coverage values at different level of percentile.

Dependencies:

660 Tie-Array-Packed-0.13: <http://search.cpan.org/~salva/Tie-Array-Packed-0.13/lib/Tie/Array/Packed.pm>

662 [9-preFiltering percentile](#):

664 **Input:**

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

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

670 2. In “~/Desktop/SeqCap/data/pre-filtering/”, generate a merged, sorted bam for
672 all samples:

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

678 “~/Desktop/SeqCap/data/pre-filtering/merge_sorted.bam” is the input bam.

680 3. A bed file:

682 ~/Desktop/SeqCap/data/in_target_assemblies/In_target
/all_assemblies_targetedRegionAndFlanking.bed” is generated by [5-FindingTargets](#)

684 **Commands:**

686 # Run [9-preFiltering percentile](#) under “~/Desktop/SeqCap/data/pre-filtering/”:

```
ke@NGS:~/Desktop/SeqCap/data/pre-filtering$ 9-preFiltering percentile -b
~/Desktop/SeqCap/data/pre-filtering/merge_sorted.bam -o CGRL -B
~/Desktop/SeqCap/data/in_target_assemblies/In_target/all_assemblies_targetedRegionAndFlanking.bed
```

692 **Output:**

694 In the folder “~/Desktop/SeqCap/data/pre-filtering/” there are a couple of files
created:

1. "CGRL_gene_outside_percentile.txt" shows a list of contigs having coverage <Xth or >Yth percentiles of the data. X and Y are defined by users. This file will be used in [10-SNPcleaner](#).

2. "CGRL_site_depth_percentile.txt" shows base coverage at different level of percentiles. The information in this file will be used by [10-SNPcleaner](#).

3. "CGRL_gene_depth_percentile.txt" shows average gene coverage at different level of percentiles.

4. "CGRL_gene_depth.txt" shows average coverage of each contig. If you want to know more about empirical coverage distribution of your data then you take the coverage value from this file and use R to plot it.

```
ke@NGS:~/Desktop/SeqCap/data/pre-filtering$ R
> per <- read.table("~/Desktop/SeqCap/data/pre-filtering/CGRL_gene_depth.txt",
head=F)
> hist(per$V2, breaks = 100, xlab = 'gene coverage', col = 'grey50')
> abline(v = 4.86, col = 'red', lwd = 2) #99th percentile = 4.86
```

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.

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.

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 distribution of the data. In reality per-base depth of exon capture data rarely follows a normal distribution.

722

724 **10-SNPcleaner**: Raw variant filtering and generates a “keep” file for the following
726 SNP/genotype calling by ANGSD. This script is mainly for filtering data at contig and
site levels. Users need to perform individual-level filtering before running this
script. See below for more details.

728

730 Before we call SNPs /genotypes and estimate allele frequencies using ANGSD, we
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
732 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
734 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
736 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
738 quality (e.g. potential cross-sample DNA contamination, sequencing errors,
paralogy).

740

1. Filtering at individual level

742 a. Remove individuals having extremely low or high coverage. Individual coverage
can be estimated using [8-ExonCaptureEvaluation](#) *Evaluation*. The file you want to
744 examine is “~/Desktop/SeqCap/data/ExonCapEval/data_metrics.txt”

746 `ke@NGS:~/Desktop/SeqCap/data$ less -S ExonCapEval/data_metrics.txt`

748

b. Remove individuals with excessively high sequencing error rates measured as the
750 percentage of mismatched bases out of the total number of aligned bases in the
mitochondrial genome. Empirical error can be estimated using [6-AssemblyEvaluation](#) *COVERAGE*
752

754 To run [6-AssemblyEvaluation](#) *COVERAGE* you need first to generate pileup files for
mitochondrial locus for each sample.

756

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

758

Usage 6-AssemblyEvaluation COVERAGE [options]

760

Options:

762 `-p DIR` *folder containing all pileup*
files generated by "samtools
764 *mpileup -f ref.fa sample1.bam*
> sample1.pileup"

766 -c INT coverage cutoff [5]
768 -q INT base quality cutoff [13]

770 2. Filtering at contig level

772 a. Remove contigs that show extremely low or high coverage based on the empirical
774 coverage distribution across all contigs. [9-preFiltering percentile](#) can be used to
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
files in [10-SNPcleaner](#) for the purpose of filtering.

776 b. Remove contigs with at least one SNP having allele frequencies highly deviating
778 from Hardy–Weinberg equilibrium expectations. Done by [10-SNPcleaner](#). Note this
is a very stringent filter even for exon capture dataset and not suitable at all for
780 genomic dataset. To use this filter you need to provide
“/home/ke/Desktop/SeqCap/data/in_target_assemblies/In_target/all_assemblies_c
782 ontig.bed” generated by [5-FindingTargets](#).

784 3. Filtering at site level

786 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%)
percentiles of base coverage you need run [9-preFiltering percentile](#) to get
788 “CGRL_site_depth_percentile.txt”.

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

794 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
796 sites that show a bias towards SNPs coming from the forward or reverse strand.
These will be done by [10-SNPcleaner](#).

798

-> Strand Bias: Tests if variant bases tend to come from one strand.

800 -> 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.

802 -> Base Quality Bias: Tests if variant bases tend to occur with a Phred-scale quality
bias.

804 -> Mapping Quality Bias: Tests if variant bases tend to occur with a mapping quality
bias.

806

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

810

e. Remove sites with a root mean square (RMS) mapping quality for SNPs across all

812 samples below a certain threshold. It is a measure of the variance of quality scores.
This will be done by [10-SNPcleaner](#).

814 f. (optional) For historic samples, characterize the pattern of base mis-incorporation
816 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](#).

818 Before running [10-SNPcleaner](#) make sure that individual-level filtering is finished.

820 **Input:**

822 1. “~/Desktop/SeqCap/data/pre-filtering/CGRL_gene_outside_percentile.txt” by [9-preFiltering](#) percentile.

824 2. “/home/ke/Desktop/SeqCap/data/in_target_assemblies/In_target/all_assemblies
826 _contig.bed” generated by [5-FindingTargets](#)

828 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
830 filtering.

832 4. Create a new folder “SNPcleaning” under “~/Desktop/SeqCap/data/” and inside
834 this folder generate a raw vcf that contains all sites from all individual samples that
pass individual-level filters:

```
ke@NGS:~/Desktop/SeqCap/data$ mkdir SNPcleaning
ke@NGS:~/Desktop/SeqCap/data$ cd SNPcleaning/
ke@NGS:~/Desktop/SeqCap/data/SNPcleaning$ samtools mpileup -D -I -S -uf
838 ~/Desktop/SeqCap/data/reference/all_assemblies_targetedRegionAndFlanking.fasta
~/Desktop/SeqCap/data/alignment/*sorted.bam | bcftools view -cg - > raw.vcf
```

840

842 -D output per-sample DP in BCF
-I do not perform indel calling
844 -S output per-sample strand bias P-value in BCF

846 **Commands:**

#Run [10-SNPcleaner](#) under “~/Desktop/SeqCap/data/SNPcleaning/”:

```
848 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-
850 filtering/CGRL_gene_outside_percentile.txt -X
~/Desktop/SeqCap/data/in_target_assemblies/In_target/all_assemblies_contig.bed -g
852 -v raw.vcf> out.vcf
```

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

856

Output:

858 In "~/Desktop/SeqCap/data/SNPcleaning/", several files are created:

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

#To generate a keep file for ANGSD:

864 *ke@NGS:~/Desktop/SeqCap/data/SNPcleaning\$ cut -f1,2 CGRL.bed > CGRL.keep*

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

870 #To view this file:

ke@NGS:~/Desktop/SeqCap/data/SNPcleaning\$ bunzip2 -c CGRL_filtered | less -S

872

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

Questions:

- 878 1. Check how many sites are present before and after filtering?
880 2. Check why some sites are filtered out by examining "CGRL_filtered".

882 *11-PopGenTools*

884 PopGenTools is a pipeline that relies on ANGSD (Analysis of Next Generation
886 Sequencing Data; <http://popgen.dk/wiki/index.php/ANGSD>) and ngsTools for data
filtering, SNP/genotype calling, allele frequency/theta estimate, and some basic
population genetic analyses.

888
890 In this workshop, we will use PopGenTools to call ANGSD in order to do a quick SNP
and genotype calling.

892 ANGSD is a software for analyzing next generation sequencing data. The software
can handle a number of different input types from mapped reads to imputed
894 genotype probabilities. Most methods take genotype uncertainty into account
instead of basing the analysis on called genotypes. This is especially useful for low
896 and medium depth data. The software is written in C++ and has been used on large
sample sizes.

898
900 ****To better demonstrate how the output files should look, we will use more
samples to run PopGenTools.

902 **Input files:**

904 1. Indexed BAM files generated by [7-Alignment](#)

These bam files are located in
906 "~/Desktop/SeqCap/denovoTargetCapture/associated_files/ANGSD/alignment/"

908 2. A .txt file that contains the path and name of each bam file. This file is made and
save in the backup folder

```
910 ke@NGS:~/Desktop/SeqCap/data$ head  
~/Desktop/SeqCap/denovoTargetCapture/associated_files/ANGSD/bamlist.txt
```

912 3. A .keep file that is generated by [10-SNPcleaner](#)

914 ~/Desktop/SeqCap/denovoTargetCapture/associated_files/ANGSD/CGRL.keep"

916 4. A reference genome (in_target_assemblies) generated by [5-FindingTargets](#)

~/Desktop/SeqCap/denovoTargetCapture/associated_files/ANGSD
918 /all_assemblies_targetedRegionAndFlanking.fasta"

920

Commands:

922 #create a folder "ANGSD" under "~/Desktop/SeqCap/data/"

```
ke@NGS:~/Desktop/SeqCap/data$ mkdir ANGSD
```

924

#copy all the input files to "~/Desktop/SeqCap/data/ANGSD" for your convenience

```
926 ke@NGS:~/Desktop/SeqCap/data$ cp  
~/Desktop/SeqCap/denovoTargetCapture/associated_files/ANGSD/bamlist.txt  
928 ~/Desktop/SeqCap/denovoTargetCapture/associated_files/ANGSD/all_assemblies_targetedRegionAndFlanking.fasta  
930 ~/Desktop/SeqCap/denovoTargetCapture/associated_files/ANGSD/CGRL.keep .
```

```
932 #cd to "ANGSD"
```

```
ke@NGS:~/Desktop/SeqCap/data$ cd ANGSD
```

```
934
```

```
936 #Under "~/Desktop/SeqCap/data/ANGSD", run PopGenTools to use ANGSD to call  
SNPs and Genotypes
```

```
938
```

```
ke@NGS:~/Desktop/SeqCap/data/ANGSD$ 11-PopGenTools ANGSD -b bamlist.txt -f  
940 CGRL.keep -r all_assemblies_targetedRegionAndFlanking.fasta -a  
all_assemblies_targetedRegionAndFlanking.fasta -q 2 -g 5 -n 8 -o first_try_folded -c 1 -  
942 s 1
```

```
944 Output:
```

```
1. Genotype call are stored in "first_try_folded.geno.gz"
```

```
946
```

```
ke@NGS:~/Desktop/SeqCap/data/ANGSD$ zcat first_try_folded.geno.gz | less
```

```
948
```

```
2. SNP calls and minor allele frequencies for the called SNPs are stored in  
950 "first_try_folded.mafs.gz"
```

```
952 ke@NGS:~/Desktop/SeqCap/data/ANGSD$ zcat first_try_folded.mafs.gz | less
```

```
954 3. Genotype likelihoods in genotype likelihood beagle input file format
```

```
956 ke@NGS:~/Desktop/SeqCap/data/ANGSD$ zcat first_try_folded.beagle.gz | less
```

column 1 (marker)

```
958 the chromosome and position
```

column 2 (allele 1)

```
960 the major allele codes as 0=A, 1=C, 2=G, 3=T
```

column 3 (allele 2)

```
962 the minor allele codes as 0=A, 1=C, 2=G, 3=T
```

column 4 (Ind0)

964 Genotype likelihood for the major/major genotype for the first individual
 column 5 (Ind0)

966 Genotype likelihood for the major/minor genotype for the first individual
 column 6 (Ind0)

968 Genotype likelihood for the minor/minor genotype for the first individual
 column 7 (Ind1)

970 Genotype likelihood for the major/major genotype for the second individual ...

972