

BIOM262: ChIP-Seq workshop

This workshop will walk you through an example of ChIP-seq analysis using HOMER (<http://homer.ucsd.edu/homer/>). HOMER was created by Chris Benner, and I love the documentation and tutorials. During the workshop, and in general, it is always to type the command and get the notes and use options of the command.

This workshop starts with aligned SAM files (only reads from mouse chr17) and performs many of the basic analysis tasks that one might normally do when analyzing ChIP-seq data.

0. Align FASTQ reads using bwa, bowtie2 or similar genome alignment algorithm. This will produce a SAM or BAM file that can be analyzed using HOMER. HOMER should be already installed and the mm9 genome is loaded in your environment.

1. Download the zip file containing SAM alignment files and unzip the archive.

```
cp /projects/ps-yeolab/biom262_2017/chip_seq_goren/samfiles.zip  
unzip samfiles.zip
```

The archive should contain the following SAM files that have already been aligned to the mouse **mm9** genome:

- h3k27ac-esc.chr17.2m.sam
- h3k4me2-esc.chr17.2m.sam
- input-esc.chr17.2m.sam
- klf4-esc.chr17.2m.sam
- oct4-esc.chr17.2m.sam
- sox2-esc.chr17.2m.sam

These files are originally from the following study investigating the roles that reprogramming factors play when transforming MEF (fibroblasts) into embryonic stem cells.

[Chronis et al. Cooperative Binding of Transcription Factors Orchestrates Reprogramming Sequencing Data: GSE90893](#)

For this tutorial we extracted the ChIP-seq experiments for several transcription factors and histone modifications performed ESC (embryonic stem cells). To reduce runtimes, only reads that mapped to chr17 (and chr17_random) are included in the SAM files (how would you generate such a file?).

It is a good practice to always double check datasets before you start analyzing them. For instance, use samtools to view the files (*samtools view -h data/sams/h3k27ac-esc.chr17.2m.sam | less*) and validate that the files are indeed what they should be (e.g., aligned to chr 17, and have 2M reads). To calculate the number of reads do *samtools view -h data/sams/h3k27ac-esc.chr17.2m.sam | wc -l*. If you want to understand better the way SAM files are organized you can follow <https://samtools.github.io/hts-specs/SAMv1.pdf> section 1.4.

2. Create a “tag directory” for the example Oct4 ChIP-seq experiment using the makeTagDirectory command. Start by typing *makeTagDirectory* (without any options) in your command line, it will provide the usage, some info about the command and a full list of program options – I highly recommend doing that whenever you use a new command.

Tag directories are analogous to sorted bam files and are the starting point for most HOMER operations like finding peaks, creating visualization files, or calculating read densities. The

command also performs several quality control and parameter estimation calculations. The command has the following form:

makeTagDirectory <Output Tag Directory> [options] <input SAM file1> [input SAM file2] ...

To create a tag directory for the Oct4 experiment, run the following command with recommended options:

makeTagDirectory <path to where your tag directories will be>/oct4-esc -genome mm9 -checkGC <path to your sam files>/oct4-esc.chr17.2m.sam

The command will take several seconds to run. What it is doing is parsing through the SAM file, removing reads that do not align to a unique position in the genome, separating reads by chromosome and sorting them by position, calculating how often reads appear in the same position to estimate the clonality (*i.e.* PCR duplication), calculating the relative distribution of reads relative to one another to estimate the ChIP-fragment length, calculating sequence properties and GC-content of the reads, and performs a simple enrichment calculation to check if the experiment looks like a ChIP-seq experiment (vs. an RNA-seq experiment).

The command creates a new directory, in this case named “oct4-esc”. Inside the directory are several text files that contain various QC results. Try opening the following (either on the command line by typing *less -S <filename>*, or if you prefer by loading to R.

tagInfo.txt - summary information from the experiment, including read totals.

tagFreqUniq.txt - nucleotide frequencies relative to the 5' end of the sequencing reads.

genomeGCcontent.txt - distribution of ChIP-fragment GC%

tagAutocorrelation.txt - relative distribution of reads found on the same strand vs. different strands.

tagCountDistribution.txt - number of reads appearing at the same positions.

3. Now we will create tag directories for all samples, by following using a shell ‘for loop’. First,

for f in <path to your sam files>/*.sam; do fname=`basename \$f .chr17.2m.sam`; makeTagDirectory <path to where your tag directories will be>/\$fname -genome mm9 -checkGC \$f; done

At this point you should have 7 tag directories. Look through the QC stats of the various ones – each dataset was created by a different antibody, and they can be divided into three types: TFs, HMs and global input. Since we will need to treat each type differently, I recommend making a directory for each – input, tfs and hms and move the tag directories to the relevant one (e.g. *tfs/oct4-esc/*, etc).

4. Next we will visualize the ChIP-seq experiments by creating bedGraph files from the tag directories and using the IGV genome browser to look at the results. We will do this using the **makeUCSCfile** command. For most ChIP-seq experiments all you need to do is specify the tag directory and specify “-o auto” for the command to automatically save the bedGraph file inside the tag directory:

makeUCSCfile <Tag Directory> -o auto

For a specific dataset, *e.g.* Oct4, the command would be:

makeUCSCfile <path to your tag directories>/oct4-esc/ -o auto

This creates the file “oct4-esc/oct4-esc.ucsc.bedGraph.gz”. This file format specifies the normalized read depth at variable intervals along the genome (use *zmore* and the filename to view the file format for yourself).

Now make these for all samples:

```
for dir in <path>*esc; do makeUCSCfile $dir -o auto; done
```

To view the file in the genome browser, do the following:

- Download the files to your computer (*scp ucsd-train<your number>@tsccl-login.sdsc.edu:/home/ucsd-train<your number>/<full path to the file> <path to location to be copied to>*; for instance in my environment it is:

```
scp ucsd-train36@tsccl-login.sdsc.edu:/home/ucsd-train36/data/tfs/oct4-esc/oct4-esc.ucsc.bedGraph.gz <path>/BIOM262/bedGraphs/
```

or for all files:

```
scp -r ucsd-train36@tsccl-login.sdsc.edu:/home/ucsd-train36/data/*/*/*.bedGraph.gz <path>/Teaching/BIOM262/bedGraph
```

- Open IGV. Make sure you use the right genome (mm9) and drag the file to the center window (or select file -> load from file).
- The read pileups will display the relative density of ChIP-seq reads at each position in the genome. We only have data for chr17 in this example, so stick to that chromosome.

5. See if there are any interesting patterns in the data that catch your eye. Try visiting the Pou5f1 locus (the gene for Oct4) by typing the gene name into the search bar at the top. Once at the Pou5f1 locus, zoom out (alt+click or scale on top right) to see if there any nearby sites that might resemble enhancers.

--

2/23/17

6. One of the most common tasks with ChIP-seq data is to find ‘enriched’ regions commonly called “peaks”. HOMER contains a command called findPeaks which is used to analyze tag directories for peaks. There are two common ways to use the command:

```
findPeaks <tag directory> -i <control tag directory> -style factor -o auto  
findPeaks <tag directory> -i <control tag directory> -style histone -o auto
```

The difference between the two is in the “-style factor/histone” argument, which will tell the program to look for focal, fixed width peaks vs. variable length peaks; the later is more common in the case of histone modifications. To find Oct4 peaks in the data, run the following command:

```
findPeaks <path>/oct4-esc/ -i <path>/input-esc/ -style factor -o auto
```

This command will look for enriched regions and filter them based on several criterion, including ensuring that they have at least 4-fold more reads in peak regions relative to the control experiment (in this case “input-esc/”). The output will be stored in a HOMER-style peak file

located in the Oct4 tag directory (“oct4-esc/peaks.txt”). The beginning of this file contains statistics and QC stats from the peak finding, including the number of peaks, number of peaks lost to input filtering, etc.

One field worth paying attention to is the “**Approximate IP efficiency**” which reports what fraction of reads from the experiment were actually found in peaks. For most decent experiments this value ranges from 1% to >30% (remember ChIP is an enrichment strategy... there is plenty of background in the data too!). Below this are the peaks along with enrichment statistics for each region.

One other thing to note is that HOMER reports the results in a ‘peak’ file, which has a slightly different format from a traditional BED file format. To create a BED file from the peak file, use the tool `pos2bed.pl` (i.e. **pos2bed.pl oct4-esc/peaks.txt > oct4-esc.bed**). BED files can be uploaded to IGV just like a bedGraph file. Also, most HOMER programs will work with either BED or peak files as input.

Next we will find peaks for all samples using two ‘for loops’ – for the two types of data:

```
for dir in <path>/hms/*; do findPeaks $dir -i <path>input-esc/ -style histone -o auto; done  
and  
for dir in <path>/tfs/*; do findPeaks $dir -i <path>input-esc/ -style factor -o auto; done
```

7. Now that we have identified peaks from our ChIP-seq data, it is time to figure out more information about where they are and what genes they might be regulating. HOMER contains a program called `annotatePeaks.pl` that performs a wide variety of functions using peak/BED files. First, let's use it to perform basic annotation of the peak file. The `annotatePeaks.pl` program works like this:

```
annotatePeaks.pl <peak/BED file> <genome version> [options] > output.txt
```

The “> output.txt” part at the end means that the results will be sent to stdout, and the “> output.txt” is used to capture the output information in a file. To annotate peaks from the Oct4 experiment:

```
annotatePeaks.pl <path>/oct4-esc/peaks.txt mm9 > oct4.annotation.txt
```

If we view the “oct4.annotation.txt” file with `less -S`, you'll see several annotation columns. Take note of the columns specifying the nearest gene TSS, the distance, and the annotation of the genomic region the peak is located in. This annotation is split into two separate columns - one is basic (i.e. exon, promoter, intergenic, intron etc.), and a more detailed annotation that describes CpG islands, repeat elements, etc. You might have also noticed while the command was running that stats about annotation enrichment too.

You can do it using a for loop for the two types of datasets:

```
for dir in hms/*; do dirname=${dir##*/}; annotatePeaks.pl $dir/regions.txt mm9 >  
annotations/$dirname.annotation.txt; done
```

and for tfs:

```
for dir in tfs/*; do dirname=${dir##*/}; annotatePeaks.pl $dir/peaks.txt mm9 >  
annotations/$dirname.annotation.txt; done
```

8. The `annotatePeaks.pl` program can also be used to create histograms that display the relative read enrichment relative to given genomic features, including transcription start sites (TSS) or

any other set of regions the user wants to define. Since the TSS is so commonly used for this purpose, HOMER has a built-in annotation for TSS (based on RefSeq transcripts). The key parameters to create a histogram are the “-hist #” and “-size #” options, which control the binning size and total length of the histogram. The other important option is the “-d <tag directory>”, which specifies which experiments to compile histograms for. In general:

```
annotatePeaks.pl <peak/BED file> <genome version> -size <#> -hist <#> -d <Tag Directory1> > output.txt
```

(note that the peak/BED file can be replaced with the key word “tss” to make a histogram at the TSS). To create a histogram with the experiments we’ve looked at thus far near the TSS, run the following:

```
annotatePeaks.pl tss mm9 -size 8000 -hist 10 -d <path>/oct4-esc/ <path>/sox2-esc/ <path>/h3k27me3-esc/ <path>/input-esc/ > output.txt
```

Open the “output.txt” using R. You’ll notice that the first column gives the distance offsets from the TSS followed by columns corresponding to the ‘coverage’, ‘+ Tags’, and ‘- Tags’ for each experiment. Try graphing each as X-Y line graph using the first column as the X-coordinate to see the patterns.

9. DNA motif finding is a powerful technique to analyze ChIP-seq experiments. Unlike gene expression data, ChIP-seq localizes signals to very specific regions of the genome allowing for accurate identification of the genetic signals responsible for recruiting various transcription factors. To use HOMER’s motif analysis program, run the findMotifsGenome.pl command using peak files from the experiments. In general the command works like this:

```
findMotifsGenome.pl <peak/BED file> <genome version> <output directory> [options]
```

Common options for motif finding are “-p <#cpu>” for parallel execution, and “-size <#>” to specify the size of the regions you wish to search for DNA motifs. For transcription factor motifs, a good size is 100. To find Oct4 enriched motifs, run the following command:

```
findMotifsGenome.pl oct4-esc/peaks.txt mm9r motifs-oct4/ -size 100 -p 10
```

or for all:

```
for dir in <path>/forLoop/tfs/*; do dirname=${dir##*/}; findMotifsGenome.pl $dir/peaks.txt mm9r <path>/motifs-$dirname/ -size 100 -p 10; done
```

This command will perform several different steps, including checking for the enrichment of a library of known transcription factor motifs as well as perform a *de novo* search for enriched motifs. The “mm9r” tells the program to mask repeat sequences in the genome, but it will still work well if you simply specify “mm9” as well. Depending on the speed of your computer this program may take while to run. Once it’s finished, download the homerResults.html file located in the motifs-oct4/ directory. Copy to your computer, and use your internet browser to open it.

This file will list the top *de novo* motifs found as well as provide stats and best matches to known transcription factor motifs. A second file called knownResults.html contains the enrichment statistics for known motifs.

You might notice that both analyses indicate that Oct4 and Sox2 peaks in ESC are highly enriched for an OCT:SOX composite motif, which has been shown to be very important in establishing pluripotent enhancers.

10. Now that you've found the most enriched motifs in a ChIP-seq experiment, it is worth it to see where those motifs are located (i.e. which peaks, etc.). One of the key outputs from motif finding are "motif files", which contain the information needed to understand where the motif is located in the genome. For example, the top de novo motifs found during motif finding are located in the output directory in the homerResults/ directory. To make it easier, let's copy the top Oct4 motif to the file "topOct4.motif" in the main directory where we are executing these commands (alternatively you could save the motif file down from the HTML results):

```
cp motifs-oct4/homerResults/motif1.motif topOct4.motif
```

Now let's perform two separate analyses - first, let's create a histogram showing the motif positions relative to Oct4 peaks so we can check if it really looks enriched relative to the center of the Oct4 peaks. We can do this using the annotatePeaks.pl program like before, but instead of looking at read densities with the "-d <tag directory>" option we can look at motif densities using the "-m <motif file>" option instead:

```
annotatePeaks.pl oct4-esc/peaks.txt mm9 -size 2000 -hist 10 -m topOct4.motif > Oct4motifHistogram.txt
```

Once the command finishes, open **Oct4motifHistogram.txt** in R and create an X-Y plot to examine the distribution of the motif relative to the peaks.

Next, let's see where these motifs are located in the actual genome browser. To do this, we can run annotatePeaks.pl again, but this time we will not make a histogram and instead use the "-mbed <bedfile>" option to tell it to create a bed file of the motif positions that we can upload to the genome browser. Try the following:

```
annotatePeaks.pl oct4-esc/peaks.txt mm9 -mbed topOct4.motifTrack.bed -m topOct4.motif > output.txt
```

Now you can load the "**topOct4.motifTrack.bed**" file as a custom track IGV. In addition, the other output file "output.txt" will contain the peak annotation results with an additional column showing the peaks that contain the given motif.

11. Finally, we want to gain some experience comparing ChIP-seq experiments. At first pass it might make sense to make a Venn diagram comparing peaks from two experiments to see how many overlap. This is a *horrible* way to analyze ChIP-seq data!!! This is because many peaks are close to the threshold of detection, barely making the cut for statistical significance (or not) in one experiment or another. A good practice is to create a scatter plot comparing the read counts between two experiments directly at all of the sites where there is signal (i.e. peaks). To do this, first let's merge the peak files from the two experiments, collapsing peaks found that overlap:

```
mergePeaks oct4-esc/peaks.txt sox2-esc/peaks.txt > oct4andsox2.peaks.txt
```

When you run this command, you'll notice that it will print out the numbers of overlapping and unique peaks from each file (i.e. Venn diagram). This can be useful to give you a general idea of how similar the experiments are, but be careful not to over interpret these values. Now that you have the combined features in the "oct4andsox2.peaks.txt" file, let's use annotatePeaks.pl to quantify the read counts from each experiment at each of the peaks by specifying each tag directory with the -d option like the following:

```
annotatePeaks.pl oct4andsox2.peaks.txt mm9 -d oct4-esc/ sox2esc/ > scatter.txt
```

We can not open the **scatter.txt** file in R to directly look at the read counts at each peak to find those with low levels in one experiment or the other. Try creating an X-Y scatter plot of the read counts with log-transformed axes to get a sense for how different (or similar) each experiment is from one another. Now that we have an appreciation for how similar the experiments are, lets try using the `getDifferentialPeaks` command to selective find peaks that are 'specific' to one experiment relative to another. In the following example we'll specify "-F 2" to indicate that we want to find Sox2 peaks that are 2-fold higher in the Sox2 experiment relative to the Oct4 experiment:

```
getDifferentialPeaks <peak file to check> <target tag directory> <background tag directory> -F <fold change> > output.txt
```

```
getDifferentialPeaks forLoop/tfs/sox2-esc/peaks.txt forLoop/tfs/sox2-esc/ forLoop/tfs/oct4-esc/ -F 2 > sox2-specific-peaks.txt
```

Now that we have found Sox2 peaks that are relatively uniquely bound by Sox2 (and not Oct4), look at the peak file using **less -S**, and look at several of the peaks in the genome browser to convince yourself that we have found interesting peaks.

Finally, lets use motif finding to see if we can identify what is unique about the DNA sequences in the Sox2 -specific peaks that might have lead to a differential recruitment of Sox2 versus Oct4 at these sites. Run motif finding on the Sox2 specific peaks:

```
findMotifsGenome.pl sox2-specific-peaks.txt mm9r motifs-sox2-specific/ -size 100 -p 10
```

10. Notice anything different about these results relative to the results from all Sox2 peaks that we found in step 9?