

ChIP-Seq Data Analysis

ChIP-Seq is a powerful method to identify genome-wide DNA binding sites for a protein of interest. This technical note describes a simple approach to building annotated tag and count tables from ChIP-Seq data sets from the Illumina Genome Analyzer.

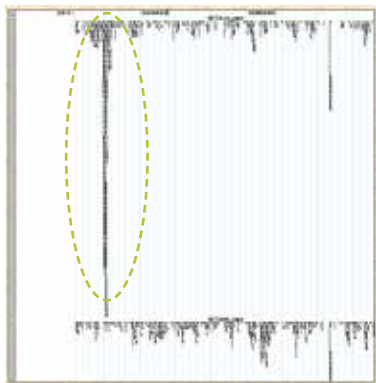
Introduction

ChIP-Seq data is less complex than other types of massively parallel sequencing data since analysis consists of determining a census count of tags from a relatively purified DNA sample. There are, however, some informatics steps that must be followed to extract meaningful data from the raw sequence reads.

The procedures described in this technical note are not intended to be complete and rigorous, but are meant to provide a starting point for researchers to successfully generate counts of sequence tags at various positions in the genome. This starting point is sufficiently flexible to facilitate novel or previously described techniques for the analysis of output data.

In addition to descriptions of how data are handled by Illumina Genome Analyzer Pipeline Software, several publicly available analysis algorithms for ChIP-Seq data analysis are discussed.

Figure 1: Chip-Seq Tags in Bed Format, Displayed in UCSC Browser



Example of custom tracks submitted in BED format (upper track is from ChIP sample and lower track is from mock control sample). The peak on the left in the ChIP sample (green circle) is significant. However, the peak on the right side is detected in both the ChIP and mock samples and is not significant.

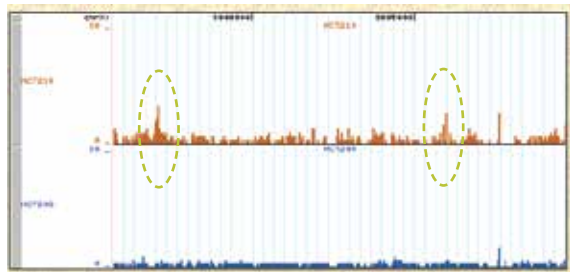
Standard Analysis Process

Illumina ChIP-Seq data produced from the Genome Analyzer are transitioned through several phases to prepare them for thorough analysis. The approach outlined below provides data display amenable for visual analysis only, which means that a researcher may have to be familiar with, or have prior expectations of a certain region of the genome to be enriched for the ChIP sequences.

1. First, the data enter the Image Analysis and Base Calling phases. Here, the actual sequence data are generated from the images acquired during sequencing by synthesis chemistry on the Genome Analyzer.
2. The short sequence reads are then aligned to the genome using ELAND. The ELAND output sequentially lists all of the sequence reads with their respective genomic coordinates. ELAND is described fully in the Genome Analyzer Pipeline Software User Guide.
3. Read data that are uniquely aligned to a genome can be viewed as a custom track in the UCSC genome browser. The track can be submitted in either a BED (Figure 1) or WIG (Figure 2) format. More information on UCSC custom tracks is available at <http://genome.ucsc.edu/goldenPath/help/customTrack.html>.
 - a. The BED format is a simple text file that contains the chromosomal start and end positions. A track in the BED format can be easily generated from GERALD *_realign.txt files. GERALD is described fully in the Genome Analyzer Pipeline Software User Guide. The following is example code that would create a BED format text file from the sequence tags in lane 3, assuming sequence reads that are 25 nucleotides in length:

```
cat s3_????_realign.txt | grep -v '^$' | \
perl -ane 'if (@F>3){$ =~/(chr.):(\
d+)\s([F|R]);print $1,"t",$2,"\
t",($2+25),"\n"}' \
> s3_customTrack.txt
```

Figure 2: Chip-Seq Tags in Wig Format



Shown is the UCSC Browser display of WIG format tracks for the same region shown in Figure 1. In this format, peaks are less obvious (green circles).

