

---

# ChiP-seq QC Report For

*H3K<sub>4</sub>me<sub>3</sub> S3*

---

qing

June 22, 2016

# ChiLin Summary

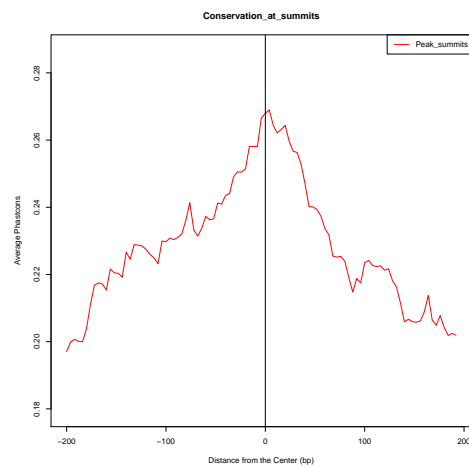
<b>I</b>	<b>Summary Table</b>	<b>2</b>
<b>II</b>	<b>FastQC</b>	<b>4</b>
II.1	FastQC score distribution . . . . .	5
II.2	FastQC GC content distribution . . . . .	6
<b>III</b>	<b>Basic mapping QC statistics</b>	<b>7</b>
<b>IV</b>	<b>Conservation</b>	<b>9</b>
<b>V</b>	<b>Motif QC analysis</b>	<b>11</b>
<b>VI</b>	<b>Library contamination</b>	<b>13</b>
<b>VII</b>	<b>FRiP</b>	<b>15</b>
<b>VIII</b>	<b>QC result instruction</b>	<b>17</b>

# Part I

## Summary Table

Metrics	treat1	treat2
FastQC	28	29
Original total reads	19.7M	20.0M
Unique mapped reads	13.3M (67.51%)	14.4M (72.2%)
Unique locations of 4M reads	3.8M (95.74%)	3.7M (93.28%)
Locations with only 1 read from 4M reads, number (ratio)	3.7M (92.18%)	3.5M (88.44%)
PBC of 4M reads	96.3%	94.8%
Fragment size of 4M reads	131	102
DHS/Promoter/Exon ratio of 4M reads	49.29%/30.33%/7.45%	54.92%/34.52%/7.25%
FRiP of 4M non-chrM reads	44.79%	53.12%
Replicates total peaks	25295	31407
Replicates 10 fold confident peaks	10458	11475
Replicates 20 fold confident peaks	1320	1041
Replicates reads correlation/replicates peaks overlap	Correlation 0.93/Overlap 24352	
Merged total/10 fold/20 fold peaks	33228/10566/806	
Top peaks overlap with union DHS number (ratio)	4994 (99.88%)	
Exon/Intron/Intergenic/Promoter ratio of peak summits	18.99%/41.05%/20.49%/19.47%	

Top peaks conservation plot



Top peaks motif analysis

denovo -20.2

# Part II

## FastQC

## II.1 FastQC score distribution

We draw the cumulative percentage plot of the FastQC sequence quality scores of all historic data and show how your new data compare(Figure 1).

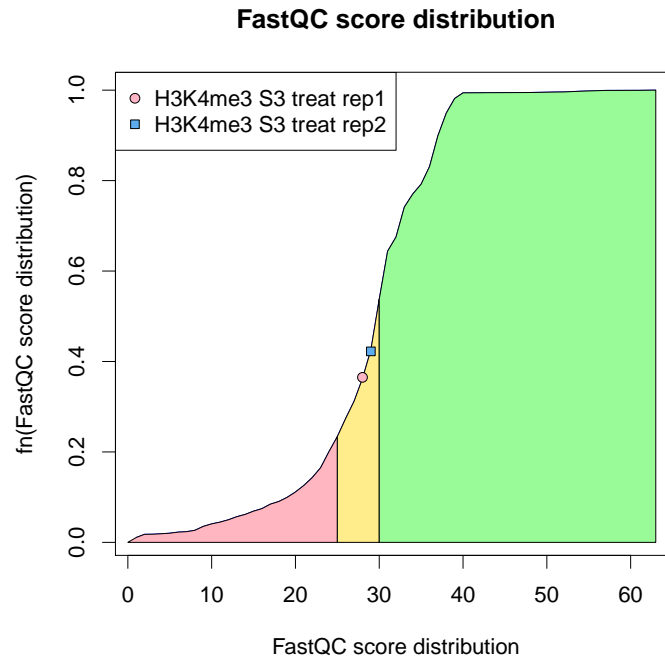


Figure 1: FastQC score distribution plot

## II.2 FastQC GC content distribution

We draw the GC content of each of replicates and controls. The expected GC content distribution (i.e. theoretical) is unimodal (Figure: 2). A bi-modal distribution may indicate sample contamination.

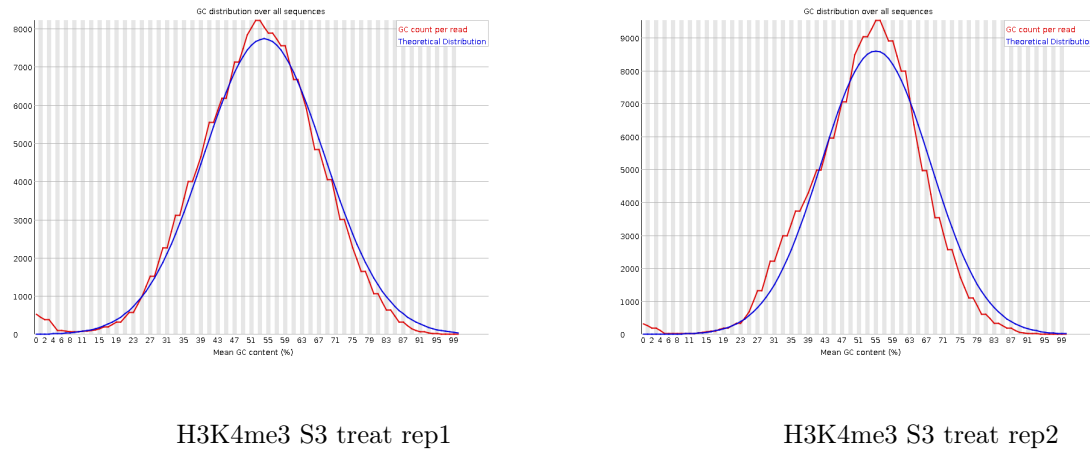


Figure 2: sequence GC content

## Part III

# Basic mapping QC statistics



We draw the cumulative percentage plot of the mappability rate of all historic data and show how your new data compare.(Figure: 3). The mappability rate is defined by the number of aligned reads divided by the number of total reads. In this plot, the x axis is the mappable rates and the y axis represents the percentage of this mappable rate among all the data. Three color range represent good, medium, and bad quality region.

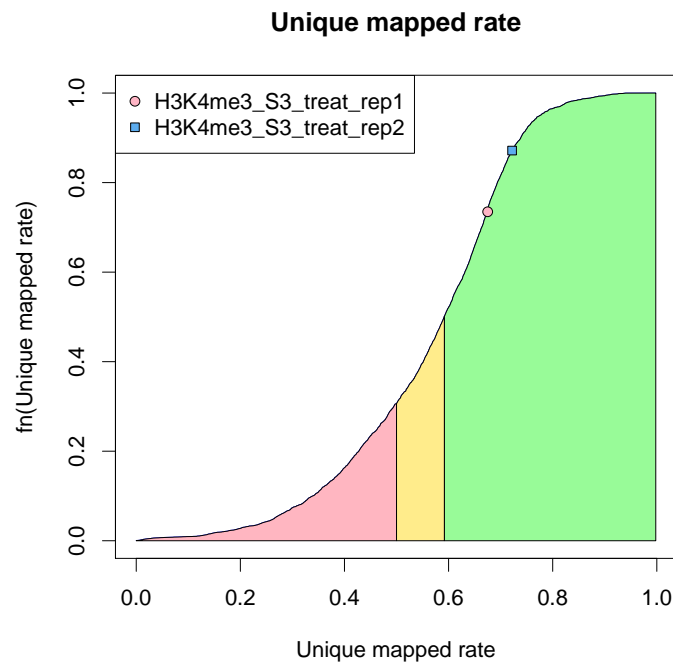


Figure 3: Mappable reads ratio

Part IV

Conservation

This section shows the conservation clustering results. Conserved sequences are similar or identical sequences across species, which have been maintained by evolution. Therefore, highly conserved sequences are thought to relate to functions conserved by evolution. We plot the average sequence conservation score of your factor binding sites within an window (TF: 500bp up and downstream; Histone: 2000bp up and downstream). The blue-plot shows the average sequence conservation score over the same window for a datasets which we chose from our historical data as “matches” to your data (i.e. “compared group”). (Figure 4)

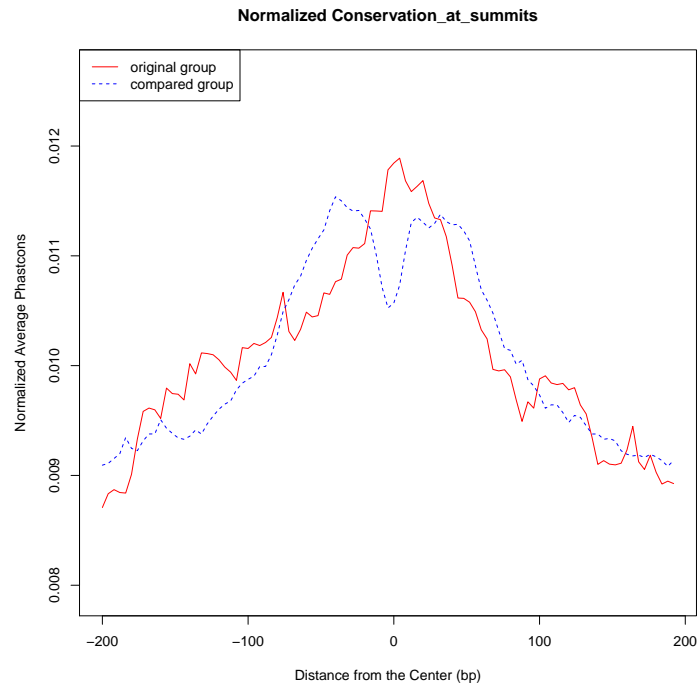


Figure 4: Phascon conservation plot

## Part V

# Motif QC analysis

Sequence motifs indicate sequence-specific binding sites for proteins such as nucleases and transcription factors (TF). <sup>1</sup> ChiLin summarize the top motif results (Table 1):

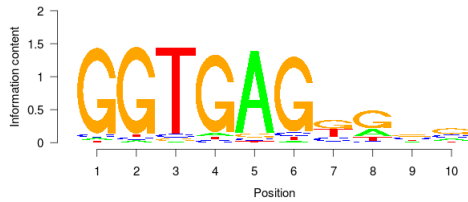
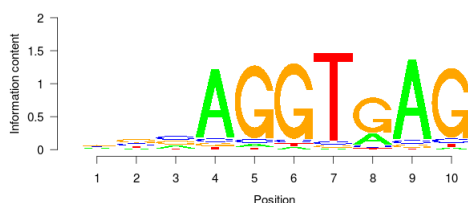
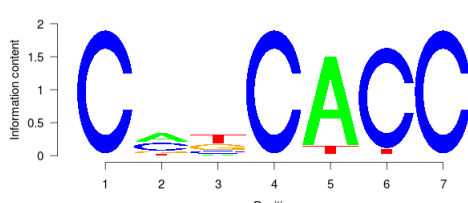

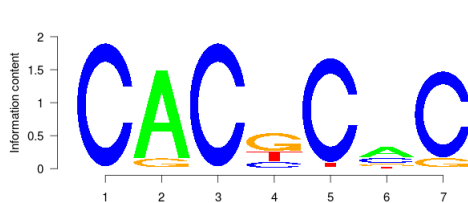
Factor name	Z-score	Hits	Motif logo
*denovo	-20.1787658763	4971	
*denovo	-17.9091841819	4674	
*denovo	-16.129750404	3814	
*denovo	-15.3247311793	1630	
*denovo	-15.2283315968	1301	

Table 1: Seqpos QC measurement

<sup>1</sup> Patrik D'haeseleer. What are DNA sequence motifs. nature biotechnology 2006,24.

## Part VI

# Library contamination

To check for cross-species contamination of the sample, we randomly select 100K reads from your files, and try to map it to the species defined in the contamination module (in the configuration file). We report the mapping rates (of the random sub-sample) for the various species.(Table 2)

sample name	mycoplasma
H3K4me3 S3 treat rep2	0.02%
H3K4me3 S3 treat rep1	0.01%

Table 2: Library contamination

## Part VII

### FRiP



FRiP (Fraction of Reads in Peaks), is a quality measure for ChIP-Seq and Dnase-Seq experiments. FRiP is simply the percentage of all tags that fall in 'hotspots', regions of local enrichment of short-read sequence tags mapped to the genome, i.e. MACS peaks. FRiP scores are positively correlated with true positive rates and replicate reproducibility.

Sample	Reads under Peaks	Total Reads	FRiP
H3K4me3 S3 treat rep2	2124916	4000000	53.12%
H3K4me3 S3 treat rep1	1791429	4000000	44.79%

Table 3: FRiP measures

## Part VIII

# QC result instruction

1. **FastQC** is the sample's median sequence quality scores. ChiLin calculates these scores using the FastQC software<sup>2</sup>. A good sequence quality score is  $\geq 25$ .
2. **Original total reads** is the sample's raw total reads number.
3. **Uniquely mapped reads** is the number of reads with mapping quality above 1. First, ChiLin aligns reads onto user-specified genomes. Then, it filters the SAM files. The uniquely mapped **RATIO** is the uniquely mapped reads divided by the total reads. A good uniquely mapped ratio is  $\geq 60\%$ .
4. **Unique locations of 4M reads** is the number of genomic locations with one or more uniquely mapped reads (unique locations) from sub-sampled 4M reads. Unique locations ratio unique locations number divided by total number of uniquely mapped reads. ChiLin estimates NRF by dividing the number of unique locations by 4M sampled uniquely mapped reads. If reads are less than 4M, then ChiLin uses the total reads instead. ChiLin reports number of unique locations and the unique locations ratio. A good unique locations of 4M reads should be  $\geq 70\%$ .
5. **Locations with only 1 read from 4M reads number (ratio)** is the number of locations with read number equal to 1 (N1). The ratio is N1 divided by 4M reads unless the total reads is less than 4M, in which case the total reads is used. A good score for this metric is  $> 70\%$ .
6. **PBC of 4M reads** is N1 (see 5) divided by unique locations (see 4). A good PBC score is  $\geq 80\%$ .
7. **Fragment size of 4M reads** is in silico estimation of your size selection through maximum cross correlation. The estimation should to be close to the size selected in your experiment.
8. **Exon/DHS/Promoter ratio of 4M reads** is the estimated ratio of reads falling in these regions (from a 4M reads sub-sample). Exons regions are defined as the merged exons regions from the RefSeq gene table. Promoter regions are defined as the RefSeq TSS  $\pm 2$ kb regions. Union DHS regions are called from ENCODE II UW DNase-seq Hypersensitive regions. The IP group samples should have higher reads ratios than the control group samples.
9. **FRiP<sup>3</sup> of 4M non-chrM reads** is used for evaluating the signal to noise ratio. First, ChiLin removes chrM reads from the total reads. Then ChiLin sub-samples 4M of these reads. Finally, it calculates the ratio of the sub-sample which fall under the called peaks. A good FRiP score is  $\geq 1\%$ .
10. **Replicates total peaks** are the total peaks number called by MACS2 with fixed extension size and q value cutoff. A good peaks number depends on your experiment.
11. **Replicates 10 fold confident peaks** are the number of peaks called by MACS2 where the fold change is  $\geq 10$ .
12. **Replicates 20 fold confident peaks** are the number of peaks called by MACS2 where the fold change is  $\geq 20$ .
13. **Replicates reads correlation** is the whole genome reads pearson correlation for all replicates with resolution 146. A good correlation score is  $\geq 0.6$ .
14. **Replicates peaks overlap** is the replicates peaks overlapping number.
15. **Top peaks not overlap with blacklist regions ratio** is the ratio of the merged top 5000 peaks (ordered by MACS2  $-\log(qvalue)$ ) which do not overlap with blacklist region<sup>4</sup>. This is expected to be  $\geq 90\%$ .

<sup>2</sup><http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

<sup>3</sup>ChIP-seq guidelines and practices of the ENCODE and modENCODE consortia

<sup>4</sup><https://sites.google.com/site/anshulkundaje/projects/blacklists>

16. **Top peaks overlap with union DHS number (ratio)** is the ratio of the merged top 5000 peaks (ordered by MACS2  $-\log(qvalue)$ ) which overlap with union DHS regions. Union DHS regions are obtained from ENCODE II UW DNase-seq Hypersensitive regions. The union DHS regions was collected from 122 human datasets or 53 mouse datasets, we do not have union DHS of other species. Union DHS generation methods is consisted of three steps: 1.for peaks length longer than 300bp, trim macs2 peaks length to 300bp around macs2 summits, 2.if less than 300bp, preserve the original length, 3.merge the peaks overlap each other. This is expected to be  $\geq 70\%$ .
17. **Exon/Intron/Integenic/Promoter ratio of peak summits** is calculated using the summits of the merged peaks. ChiLin reports the ratio of overlap with exon, intron, intergenic, and promoter regions for these summits.  
We list the background Exon/Intron/Integenic/Promoter ratio here:

Assembly	Exon	Intron	Intergenic	Promoter
hg19	1.92%	36.39%	58.37%	3.32%
hg38	1.95%	36.30%	58.27%	3.47%
mm9	1.91%	32.28%	62.38%	3.44%
mm10	1.91%	32.48%	62.14%	3.46%

18. **Top peaks conservation plot** is the Phastcons conservation scores distribution around  $\pm 2kb$  of the top 5000 merged peak summits. Phastcons conservation scores are from placental mammals multiple alignment. For TFs and active histone mark the plot should be shown as a sharp peak in the center.
19. **Top peaks motif analysis** is the motif analysis performed on the top 5000 merged peak summits. These summits are used for discovering highly enriched motifs with MDSeqPos., **Z-score**: a statistical measure of the motif credibility. A good motif will usuauall have a z-score of less than -15. **Hits**: number of times the motif occurs in the top 1000 regions.