

ChIPQC Report

Overview

This report was generated using [ChIPOC](#)

The report provides both general and ChIP-seq specific quality metrics and diagnostic graphics to allow for the quantitative assessment of ChIP-seq quality.

The report is split into three main sections:

- **QC Summary** - Overview of results.
- **QC Results** - Full QC results and figures.
- **QC files and versions** - Files and program versions used in QC

QC Summary

Table 1. Summary of ChIP-seq filtering and quality metrics.

ID	Tissue	Factor	Condition	Replicate	Reads	Dup%	ReadL	FragL	RelCC	SSD	RiP%	RiBL%
cMYC_1	A549	cMYC		1	287462	14	28	98	0.3	1.5	6.6	1.8
cMYC_2	A549	cMYC		2	317537	4.5	28	96	0.14	1.1	2.8	0
CTCF_1	A549	CTCF		1	341055	17	28	185	1.9	2.5	31	1.3
CTCF_2	A549	CTCF		2	303856	7.3	28	185	1.7	1.4	13	0
E2F1_1	HeLa-S3	E2F1		1	223580	1	28	155	0.51	1.3	7.8	2
E2F1_2	HeLa-S3	E2F1		2	194919	0.66	28	97	0.44	1.4	5.4	2.8

Table 1 contains a summary of filtering and quality metrics generated by the ChIPQC package. Further information on these metrics, their associated figures and additional quality measures can be found within the related QC Results subsections.

A short description of **Table 1** metrics is provided below:

- **ID** - Unique sample ID.
- **Tissue/Factor/Condition** - Metadata associated to sample.
- **Replicate** - Number of replicate within sample group
- **Reads** - Number of sample reads within analysed chromosomes.
- **Dup%** - Percentage of MapQ filter passing reads marked as duplicates
- **FragLen** - Estimated fragment length by cross-coverage method
- **SSD** - SSD score (htSeqTools)
- **FragLenCC** - Cross-Coverage score at the fragment length
- **RelativeCC** - Cross-coverage score at the fragment length over Cross-coverage at the read length
- **RIP%** - Percentage of reads within peaks
- **RIBL%** - Percentage of reads within Blacklist regions

QC Results

Mapping, Filtering and Duplication rate

This section presents the mapping quality, duplication rate and distribution of reads in known genomic features.

Table 2. Number and percentage of mapped, duplicated and MapQ filter passing reads

ID	Tissue	Factor	Condition	Replicate	Unmapped	Mapped	Pass MapQ Filter and Dup	Total Dup%	Pass MapQ Filter%	Pass MapQ Filter and Dup%
----	--------	--------	-----------	-----------	----------	--------	--------------------------	------------	-------------------	---------------------------

cMYC_1	A549	cMYC	1	0	287462	26874	12	69	14
cMYC_2	A549	cMYC	2	0	317537	10114	3.9	70	4.5
CTCF_1	A549	CTCF	1	0	341055	41661	15	74	17
CTCF_2	A549	CTCF	2	0	303856	15908	6.2	72	7.3
E2F1_1	HeLa-S3	E2F1	1	0	223580	1542	1.6	68	1
E2F1_2	HeLa-S3	E2F1	2	0	194919	881	1.7	68	0.66

Table 2 shows the absolute number of total, mapped, passing MapQ filter and duplicated reads. The percent of mapped reads passing quality filter and marked as duplicates (Non-Redundant Fraction?) are also included.

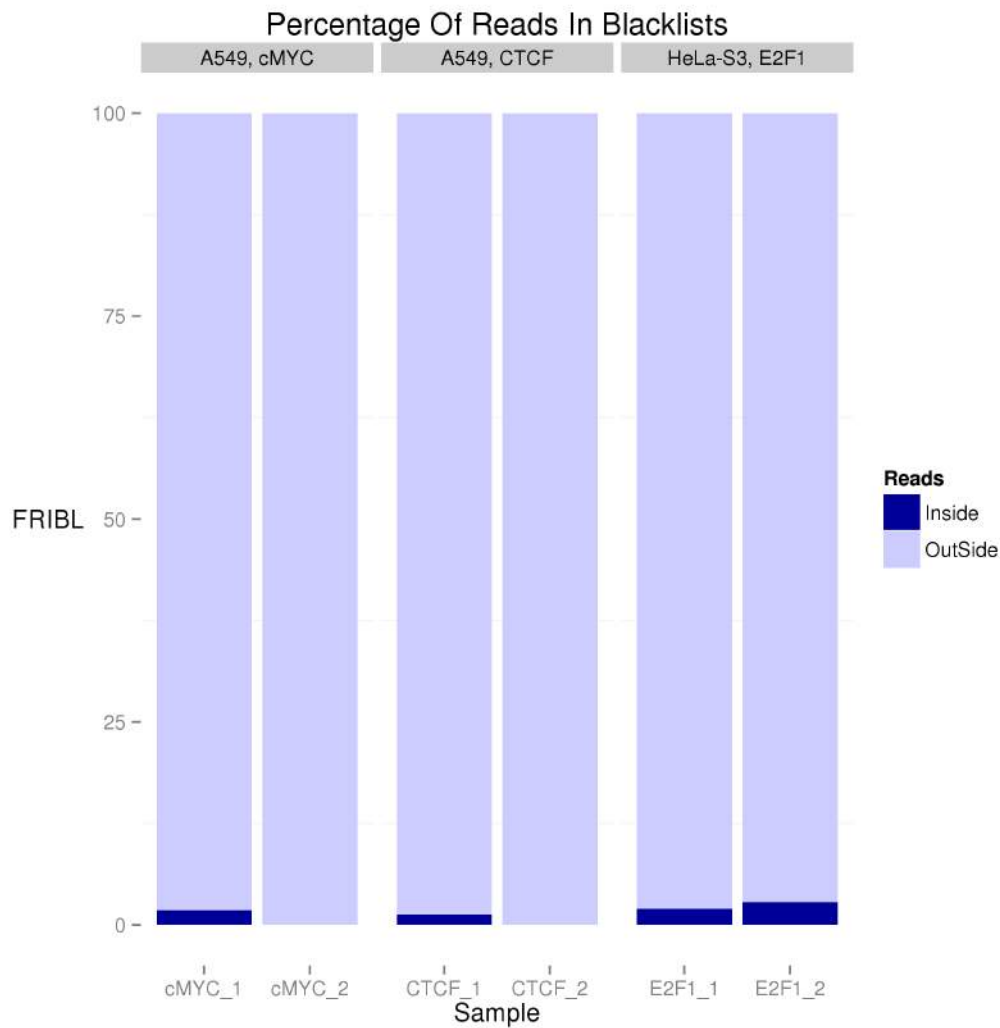
Description of read filtering and flag metrics:

- **Total Dup%**-Percentage of all **mapped** reads which are marked as **duplicates**.
- **Pass MapQ Filter%**-Percentage of all **mapped** reads which **pass MapQ quality filter**
- **Pass MapQ Filter and Dup%**-Percentage of all reads which **pass MapQ filter** and are marked as **duplicates**.

Duplication rates (Dup %) are dependent on the ChIP library complexity and the number of reads sequenced. Higher duplication rates maybe due to low ChIP efficiency when read counts are lower or conversely saturation of ChIP signal when sequencing large number of reads. Since this metric is dependent on both read depth and the properties of the ChIP itself, comparison between biological or technical replicates of similar total read counts can best identify problematic libraries .

Highly mappable (multimappable) positions within the genome can attract large levels of duplication and so assessment of duplication before and after MapQ quality filtering can identify contribution of these positions to the duplication rate.

Figure 1. Barplot of the percentage of reads in blacklists

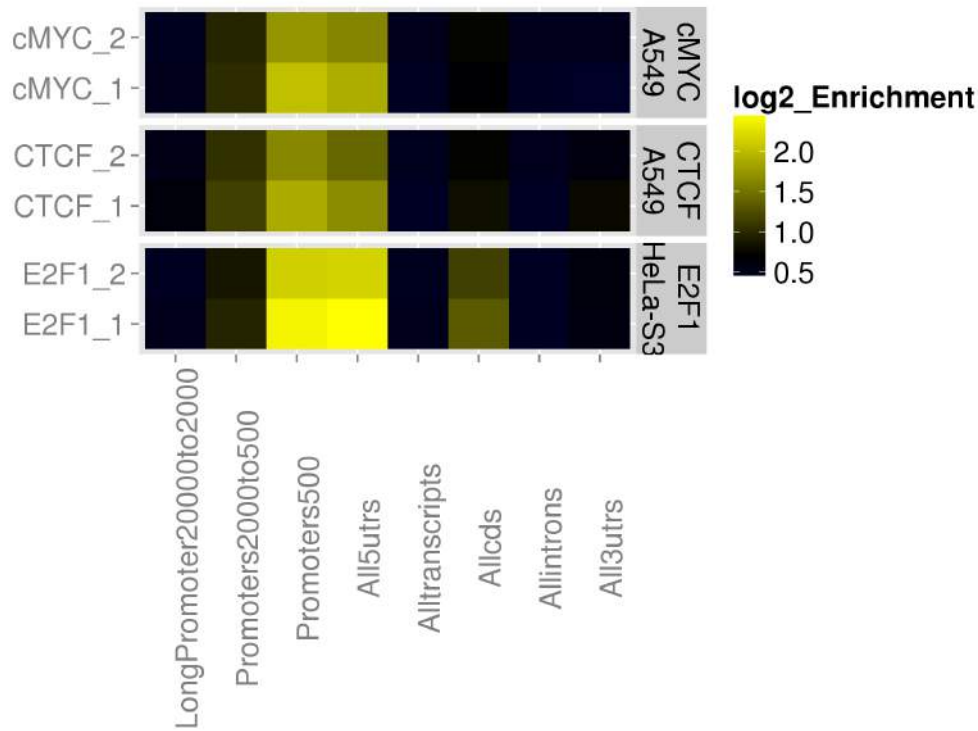


Genomic regions of high, anomalous signal have been seen to contribute directly to the Encode RCS and NSC metrics and can confound fragment length estimation, calculation of ChIP enrichment metrics (i.e. SSD) and comparison of signal between samples.

The identification of genomic stretches of artefact signal has been previously described for single samples using Input controls and more recently work as part of the Encode consortium has identified conserved regions of high artefact signal for many model organisms.

The percentage of total ChIP signal within known artefact regions can therefore be useful to evaluate the level of such confounding, abbarant signal in a sample. **(Figure 1)**

Figure 2. Heatmap of log2 enrichment of reads in genomic features



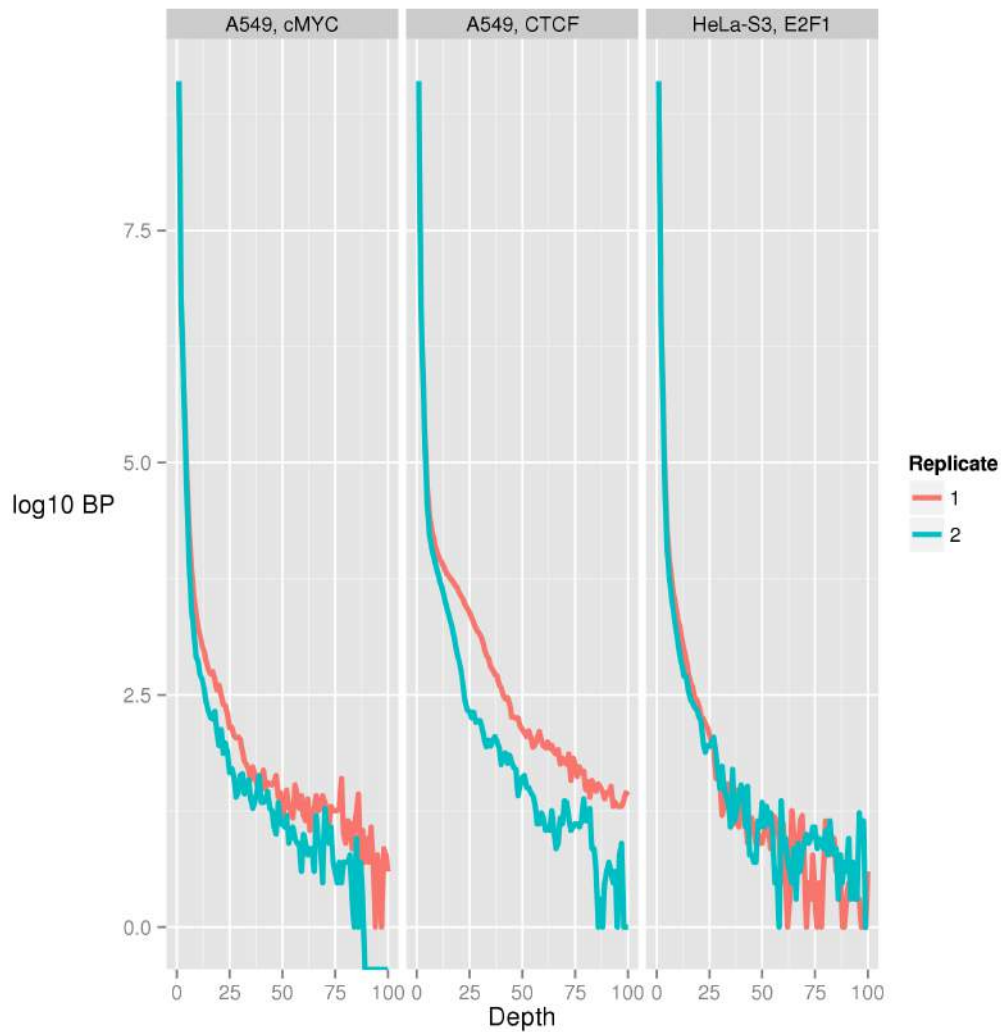
The distribution of reads across known genomic features such as genes and their subcomponents may allow further evaluation of ChIP-seq success and quality. A transcription factor known to preferentially bind at a genomic feature should show relative enrichment against other transcription factors showing no such preference. In addition, a replicate showing a differing enrichment pattern across genomic features compared to those within its sample group would highlight a potential outlier sample worthy of further investigation.

Figure 2 shows the log₂ enrichment of specified genomic features within samples with regions of greater enrichment showing bright yellow and lower enrichment seen in black.

ChIP signal Distribution and Structure

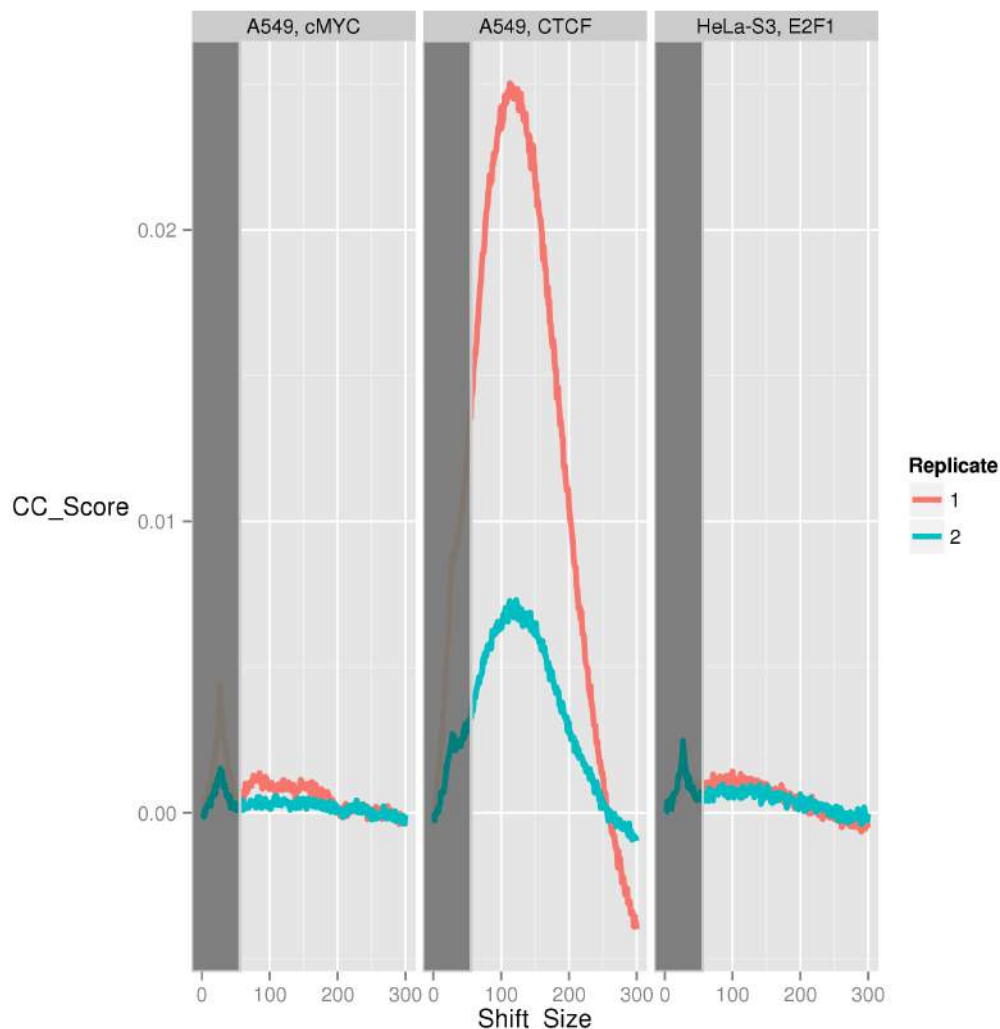
In this section, metrics relating to genome-wide depths of coverage and the relationship between Watson and Crick reads are presented. The metrics are the SSD metric and cross-coverage metrics, Relative_CC and fragmentLength_CC.

Figure 3. Plot of the log₂ base pairs of genome at differing read depths



SSD is the standard deviation of coverage normalised to the total number of reads. Evaluation of the number of bases at differing read depths, (**figure 3**) alongside the use of the SSD metric allow for an assessment of the distribution of ChIP-seq or input signal. Successful Histone and transcription factor ChIP-seq samples will show a higher proportion of genomic positions at greater depths and equivalence of sample and input SSD scores highlights either an unsuccessful ChIP or high levels of anomalous input signal

Figure 4. Plot of CrossCoverage score after successive strand shifts



An important measure of ChIP success is the degree to which Watson and Crick reads cluster around the centres of transcription factor binding sites or epigenetic marks.

Transcription factor binding sites identified by ChIP-seq will show two distinct peaks of Watson and Crick strands separated by the fragment length. Here the method of cross-coverage (ChIPseq package) analysis is used to investigate this spatial clustering of Watson and Crick reads.

To investigate this spatial clustering, reads on the positive strand are shifted in 1bp steps and the total proportion genome now covered by both strands combined is assessed.

Figure 4 shows the CCov_Score (described below) after successive shifts. The points of highest outside of the read-length exclusion region, $2 \times$ the read length, (marked in grey) is considered the fragment length

Following the methodology first presented for cross-correlation by Encode to calculate the Relative Strand Cross Correlation (NSC) and Normalised Strand Cross Correlation, the Relative Cross Coverage score and Fragment Length Cross Coverage score are calculated.

The calculation of cross-coverage (CCov), Relative CCov and Fragment Length CCov scores are explained below:

- **CCov_Score**- $1 - (\text{Total covered genome size at strand shift}) / (\text{covered genome size with no shift})$
- **Fragment Length CCov**- $(\text{CCov of fragment length strand shift}) / (\text{Minimum CCov})$
- **Relative CCov**- $(\text{CCov of fragment length strand shift}) / (\text{CCov of read length strand shift})$

Peak Profile and ChIP Enrichment

Following the identification of genome wide enrichment (peak calling), the percentage of ChIP signal within enriched regions, as well the average profile across these regions can be used to further evaluate ChIP quality

Figure 5. Plot of the average signal profile across peaks

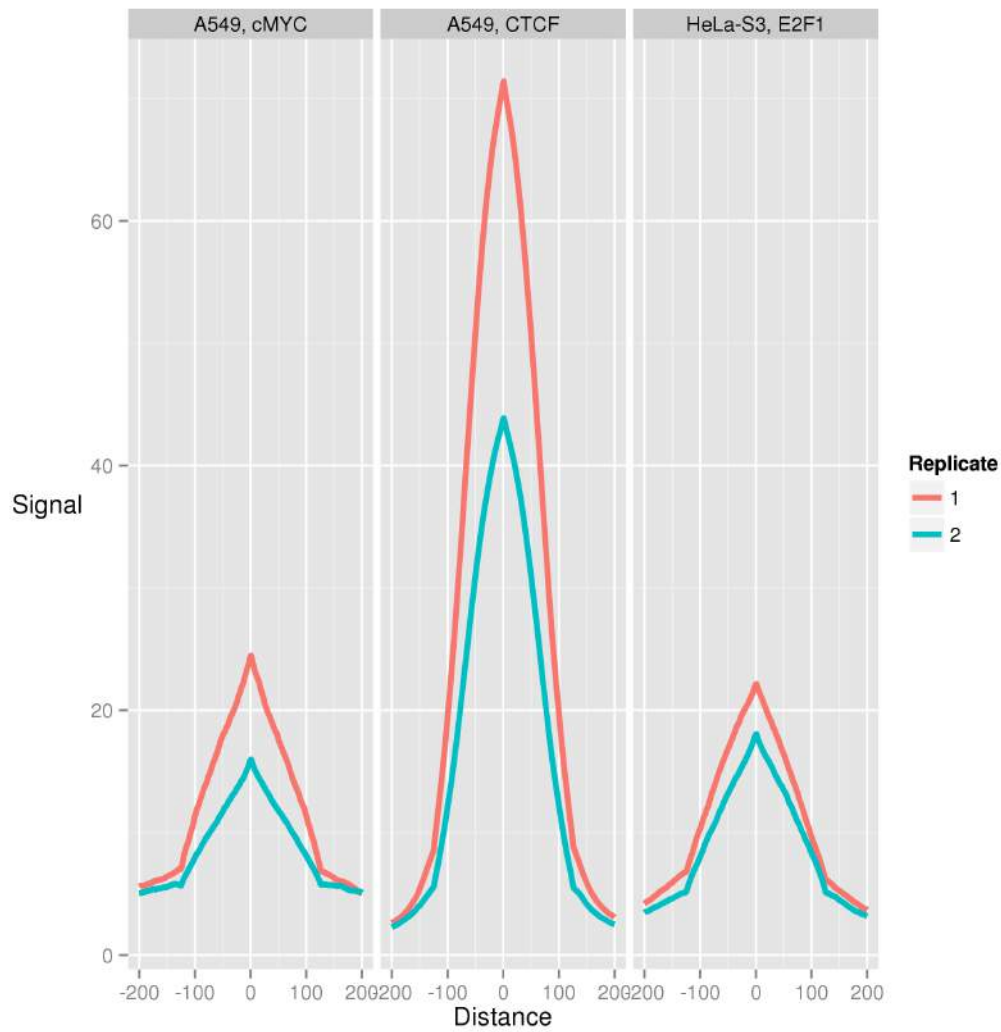


Figure 5 represents the mean read depth across and around peaks. By identifying the average pattern of enrichment across peaks, differences in both mean peak height and shape may be found. This not only assists in a better characterisation of ChIP enrichment but can aid in the identification of outliers.

Figure 6. Barplot of the percentage number of reads in peaks

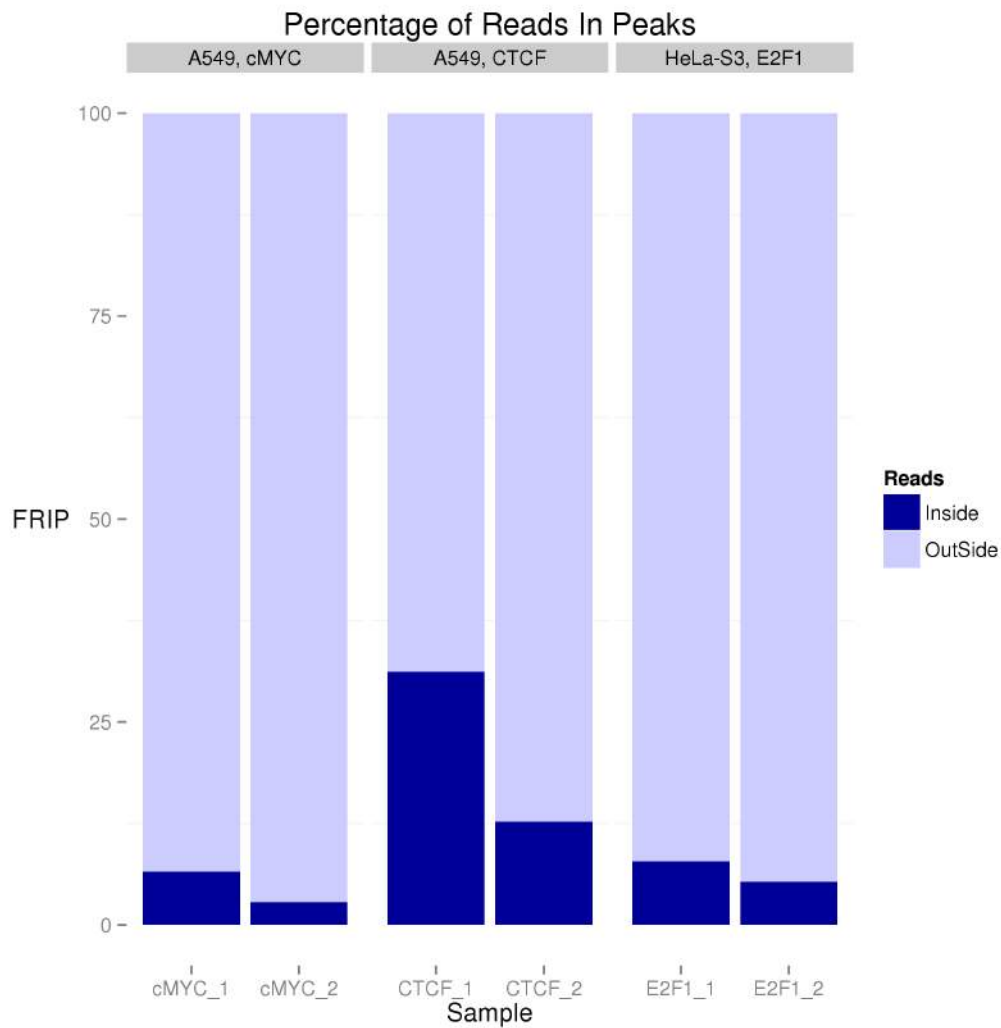


Figure6 shows the total percentage of reads contained within enriched regions or peaks. The higher efficiency ChIP-seq will show a higher percentage of reads in enriched regions/peaks and longer epigenetic marks will often have a higher ranges of efficiencies than punctate marks or transcription factors.

Figure 7. Density plot of the number of reads in peaks

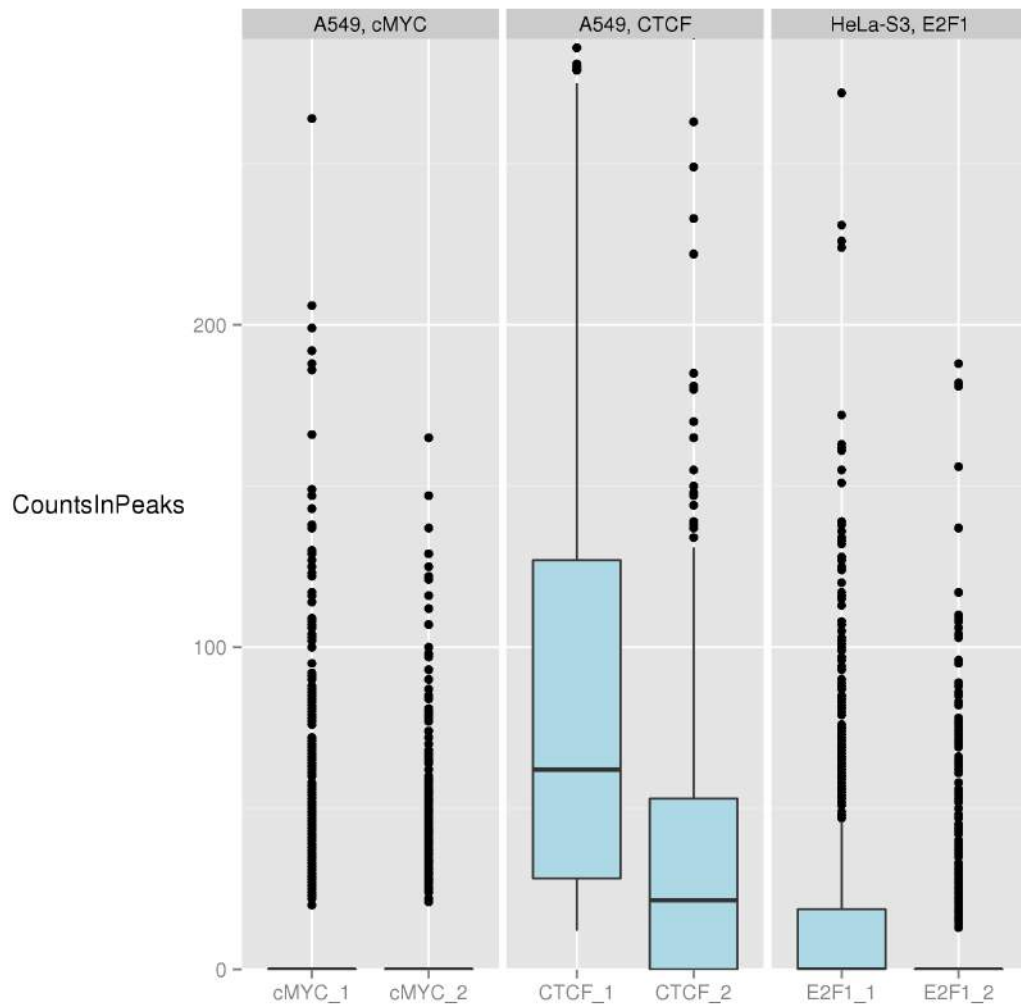


Figure 7 shows the distribution of reads in all peaks. Evaluation of the distribution can allow for greater characterisation of the variability and range of signal in peaks within a sample and so better characterise the signal across peaks than the RIP score may allow.

Figure 8. Plot of correlation between peaksets

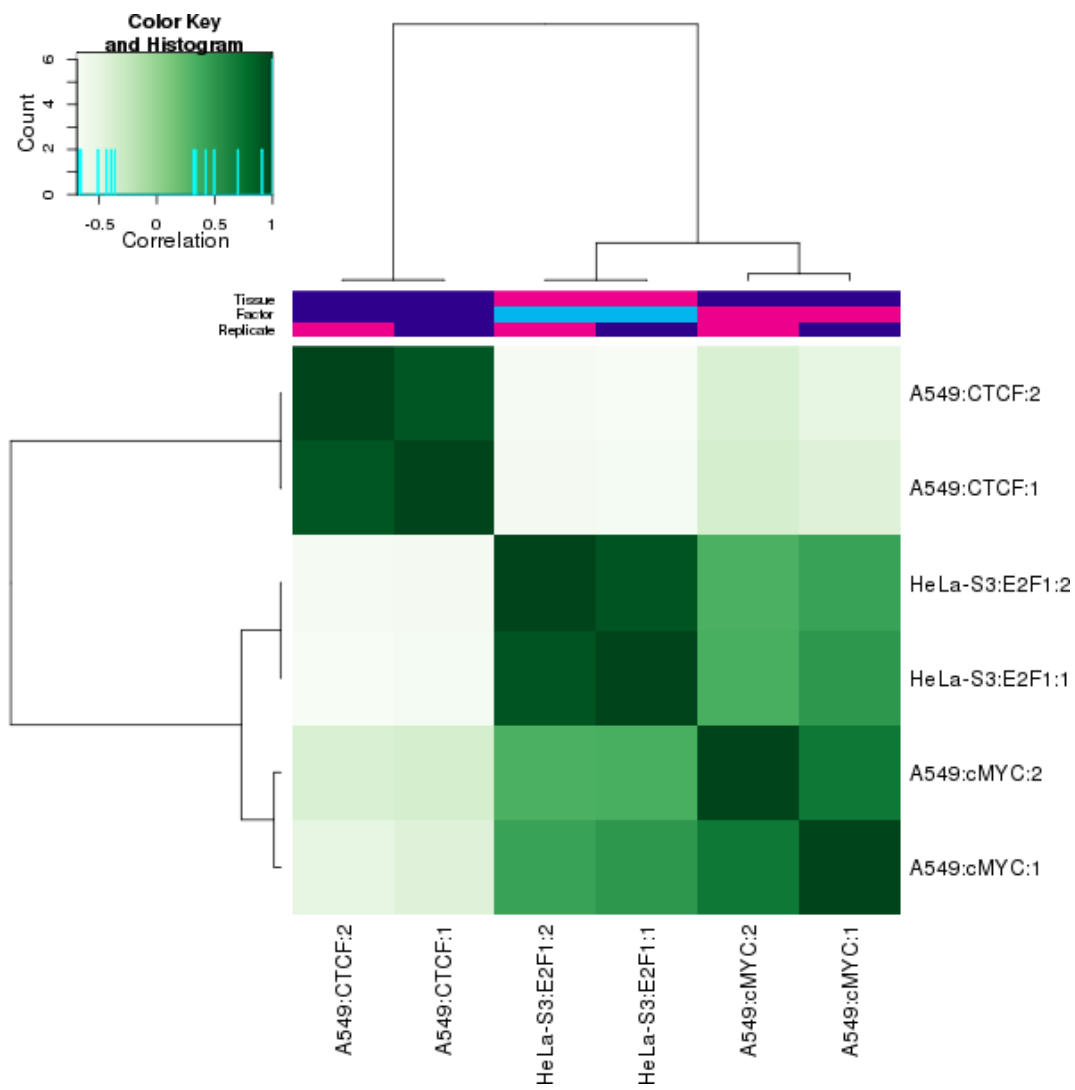


Figure 9. PCA of peaksets

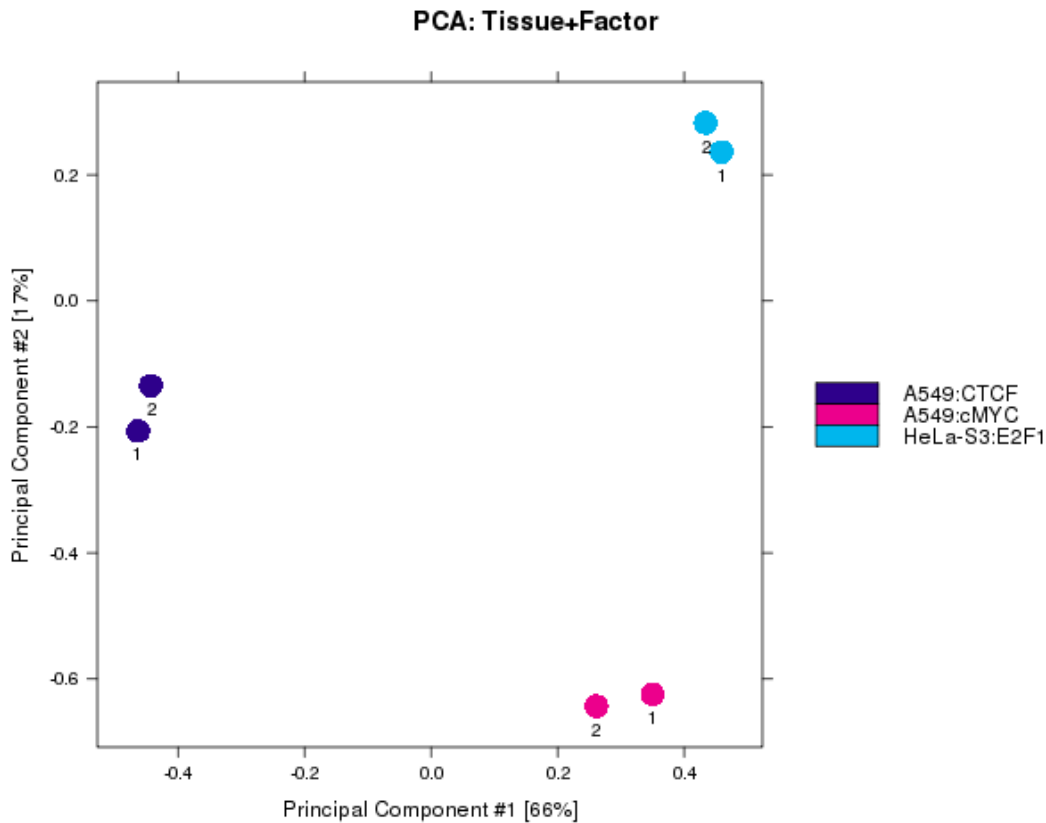


Figure8 and 9 shows the correlation between samples as a heatmap and by principal component analysis. Replicate samples of high quality can be expected to cluster together in the heatmap and be spatially grouped within the PCA plot.

Files and Versions

R Version Information

- Version: 3.1.0
- Version_String :R version 3.1.0 alpha (2014-03-18 r65213)

ChIPQC Version Information

- Version: ChIPQC:1.0.1
- Author: Tom Carroll, Wei Liu, Ines de Santiago, Rory Stark
- Maintainer: Tom Carroll , Rory Stark