

Assessing ChIP-seq sample quality with *ChIPQC*

Thomas Carroll^{1*}

¹ Medical Research Council Clinical Sciences Centre, London

*thomas.carroll (at) imperial.ac.uk

Edited: July 20th, 2014; Compiled: July 28, 2014

Contents

1	Introduction	1
2	Single sample assessment using <code>ChIPQCsample()</code>	1
2.1	A simple example of using <code>ChIPQCsample()</code>	1
2.2	Cross-coverage and the <code>FragmentLength/Relative</code> cross-coverage scores (<code>FragCC/RelCC</code>)	3
2.3	Distribution of Signal: Within peaks, blacklists and known annotation	4
2.4	Distribution of Signal: Distribution of coverage depth across the genome	6
2.5	Conclusion	8
3	Assessing a ChIP-seq experiment using <code>ChIPQC()</code>	8
3.1	An example ChIP-seq experiment using <code>ChIPQC()</code>	8
3.2	Examining Cross-coverage and <code>FragCC/RelCC</code> scores across an experiment	9
3.3	Distribution of Signal across a ChIPQC experiment	11
3.4	Assessing sample similarity with <code>Diffbind</code>	16
3.5	Conclusion	18
4	Advanced Topics	18
4.1	Providing additional data to ChIPQC plotting and reporting	18
4.2	Specifying custom annotation	20
4.3	Some plotting tips	21
5	Analysing the Tamoxifen dataset	24
6	Session Info	25

1 Introduction

This practical will cover how to assess ChIP-seq data quality using the ChIPQC package with real world datasets.

Due to limitations on time, the data has been processed and ChIPQCexperiments objects prepared for this practical. We will still cover how to set up data for ChIPQC and how to run the main processing commands but you may wish to load in the processed data provided.

2 Single sample assessment using ChIPQCsample()

2.1 A simple example of using ChIPQCsample()

ChIPQC package allows for the rapid generation of ChIP-seq quality metrics from aligned data in BAM format. The main function ChIPQCsample() can be run with simply a BAM file location and will return a ChIPQCsample object. Here we can run ChIPQCsample on a single Bam file of an EBF1 ChIP-seq experiment for chromosome 11. The data used can be found in "/data/ChIPQC/" so feel free to give the function a go or simply load the processed object.

```
library(ChIPQC)
#bamFile <- "/data/ChIPQC/Chr11_Ebf1DupMarked.bam"
#exampleExp = ChIPQCsample(bamFile, peaks=NULL, annotation=NULL, chromosomes="chr11")
load("//data/ChIPQC/ebf1_ChIPQCsample.RData")
QCmetrics(exampleExp)
```

##	Reads	Map%	Filt%	Dup%	ReadL	FragL	RelCC	SSD	RiP%
##	1.37e+06	1.00e+02	1.04e+01	6.56e+00	3.60e+01	1.77e+02	3.57e+00	1.00e+00	NA

This is the simplest way to run ChIPQCsample but to take full advantage of ChIPQC features we can add additional information on blacklisted regions, genome annotation and any peaks we have called for this ChIP.

Annotation is provided for human ("hg19", "hg18"), mouse ("mm10", "mm9"), rat ("rn4"), C Elegans ("ce6") and D Melanogaster ("dm3") by way of the Bioconductor TranscriptDb annotations packages. Additional or custom annotation can also be provided to the ChIPQCsample function and we go into this in detail later on.

Blacklist regions are included for hg19 and can be provided as a GRanges object or complete file path to blacklists in bed format. Blacklists for a range of species is available from Anshul Kundaje's google site (<https://sites.google.com/site/anshulkundaje/projects/blacklists>).

Peaks can also be provided to the ChIPQCsample function as a GRanges object or as complete file path to peaks file in bed format.

```
bamFile <- "/data/ChIPQC/Ebf1DupMarked.bam"
peaksFile <- "/data/ChIPQC/Ebf1_WithInput_Input_2_proB_peaks.bed"
```

```
BlackListFile <- "/data/ChIPQC/mm9-blacklist.bed"

#exampleExp = ChIPQCsample(bamFile, peaks=peaksFile, annotation="mm9", blacklist=BlackListFile,
load("/data/ChIPQC/ebf1_FullyLoaded_ChIPQCsample.RData")
QCmetrics(exampleExp)
```

```
##      Reads      Map%      Filt%      Dup%      ReadL      FragL      RelCC      SSD      RiP%      RiBL%
## 1.37e+06 1.00e+02 1.04e+01 6.56e+00 3.60e+01 1.77e+02 3.57e+00 1.00e+00 5.75e+00 2.23e+00
```

Now the result of `QCmetrics` contains full metrics for ChIP-seq and the additional information on `RiP%` and `RiBL%`. Lets just remind ourselves what all these mean.

Reads Total reads in bam file

Map% Percentage of total reads mapping within file.

Filt% Percentage of mapped reads passing MapQ filter.

Dup% Percentage of mapped reads marked as duplicates.

ReadL Mean read length (as integer).

FragL Predicted fragmentlength by cross-coverage method.

RelCC The relative cross-coverage score

SSD Standardised standard deviation

RiP% Reads mapped to peaks

RiBL% Reads mapped to blacklists

Now we have our full `ChIPQCsample` object we can start to review and visualise the metrics generated.

2.2 Cross-coverage and the FragmentLength/Relative cross-coverage scores (FragCC/RelCC)

For transcription factors and narrow epigenetic marks, an accumulation of Watson and Crick reads around the binding site/mark may often be seen. The degree to which your ChIP-seq signal is arranged into the Watson and Crick read clusters around such sites has been previously exploited as a metric of ChIP efficiency.

In `ChIPQC` we assess the reduction in total genome covered which occurs from shifting the Watson reads along the genome (from 5' to 3' of chromosome). This is performed by measuring measure total coverage after each every successive shift of 1bp. As the Watson reads overlap the Crick reads around peaks the total genome covered will be reduced. The total coverage after each successive shift is then converted to cross-coverage scores after each shift.

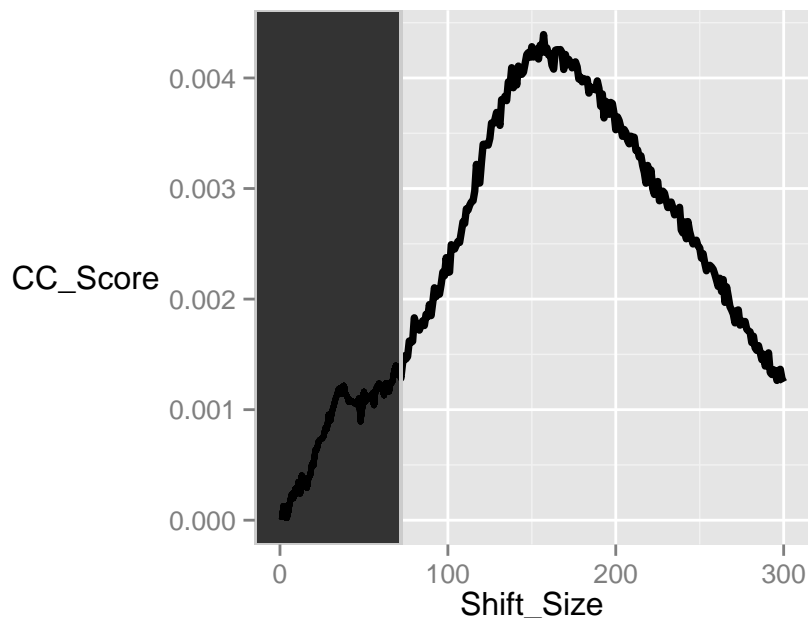
$$CrossCoverageScore_n = (TotalGenomeCoverage_0 - TotalGenomeCoverage_n) / TotalGenomeCoverage_0 \quad (1)$$

Where n is the bp shift of Watson reads and 0 is after no shift of Watson reads.

The cross-coverage scores after successive shifts can then be visualised and reviewed to identify the expected increase in cross-coverage scores around the fragment length as well as any evidence of artefacts by a peak in the cross-coverage score at the read length.

The `plotCC` function will calculate cross-coverage scores and plot those after successive shifts. Also shaded in grey is the area to be excluded when identifying the fragment length peak.

```
plotCC(exampleExp)
```



We can see that this Ebf1 ChIP has a peak in cross-coverage scores at around 160bp, corresponding to the fragment size, high above that observed at the 0 shift. This indicates that we have successfully enriched for signal around binding sites. Further to this we can also see an artefact peak at the read length and that this peak is much lower than that observed at fragment length, again indicating an enrichment for ChIP-signal over background.

Further to the visual inspection of the cross-coverage scores, we can extract the `RelCC` and `FragCC` scores for this sample using `RelativeCrossCoverage` and `FragmentLengthCrossCoverage` functions respectively.

These metrics can be considered to relate to efficiency of ChIP (`FragCC`) and efficiency of ChIP compared to artefact signal (`RelCC`) and are calculated as below.

$$FragCC = CrossCoverageScore_{max} \quad (2)$$

$$RelCC = CrossCoverageScore_{max} / CrossCoverageScore_{readlength} \quad (3)$$

Where *max* is shift with maximum cross-coverage score (excluding area 0 to 1.5*readlength) and *readlength* is the shift corresponding to the read length.

The area around the read length is excluded from selection of shift with maximum cross-coverage scores (shaded in grey in cross-coverage plots). This is due to the presence of the artefact peak which, in less enriched samples, may have greater cross-coverage score than observed at the fragment length. To avoid obscuring the selection of the fragment length peak for fragment length prediction and ascertainment of FragCC score and RelCC score this region is removed prior to calculation of shift with maximum cross-coverage score.

```
FragmentLengthCrossCoverage(exampleExp)
```

```
## [1] 0.004271
```

```
RelativeCrossCoverage(exampleExp)
```

```
## [1] 3.571
```

In this example we find, as expected from cross-coverage scores graph, that the FragCC is high and RelCC score is above 1 indicating a successful ChIP.

2.3 Distribution of Signal: Within peaks, blacklists and known annotation

Another useful set of characteristics of your ChIP-seq data is where in the genome the signal is distributed. This can be done by looking for the proportion or percentage of signal in peaks, blacklists or even in known annotation.

To get the percentage of reads landing in peaks and blacklists we can use the `rip` and `ribl` functions as well as the `plotFrip` and `plotFribl` functions for visualisation.

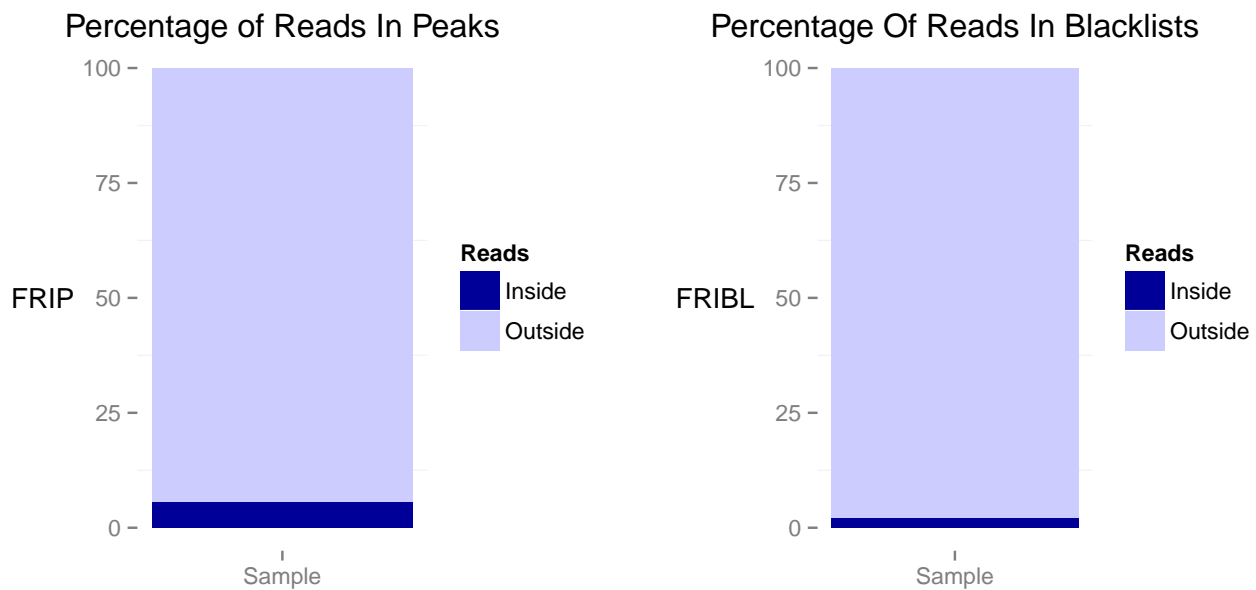
```
frip(exampleExp)
```

```
## [1] 0.05747
```

```
p1 <- plotFrip(exampleExp)
```

```
p2 <- plotFribl(exampleExp)
```

```
## Loading required package: grid
```



The frip and fribl plot, show that we have a percentage of signal in peaks greater than 5% indicating a ChIP of acceptable quality. Reassuringly, we have a higher signal in peaks than blacklisted regions (2%) and so an enrichment over artefact signal.

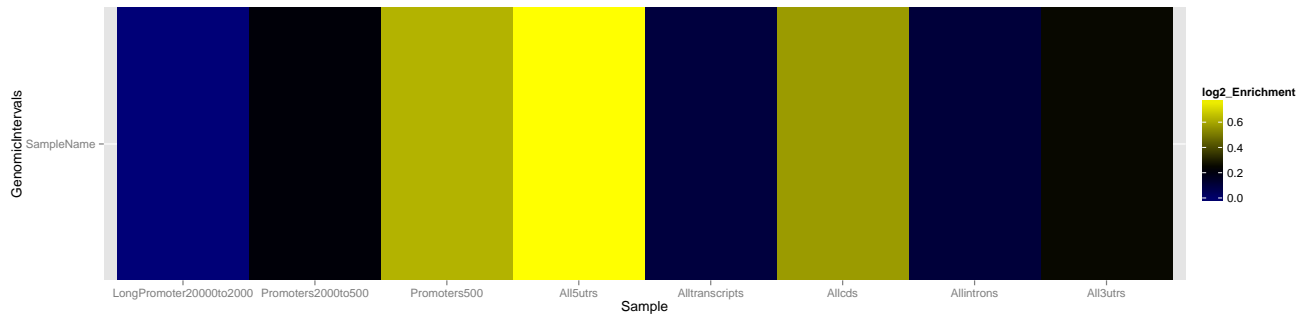
When provided with annotation in the form of genomic regions, ChIPQC can measure the enrichment of signal within them. The regi function provides a simple enrichment statistic which illustrates the distribution of signal within genomic interval annotation over that expected given their size.

$$regi = ProportionOfReadsInInterval / ProportionOfGenomeInInterval \quad (4)$$

To review regi statistics we can use the regi function or plot enrichment using plotRegi.

```
regi(exampleExp)
## LongPromoter20000to2000      Promoters2000to500      Promoters500
##           -0.02709           0.21537           0.62421
##           All5utrs           Alltranscripts           Allcds
##           0.78965           0.09785           0.57157
##           Allintrons           All3utrs
##           0.10562           0.25087

plotRegi(exampleExp)
```



The regi scores and heatmap show an enrichment for regions around the TSS including 5'UTRs, CDS (potentially first exon) and 500upstream regions. This suggests that Ebf1 is a promoter associated transcription factor as expected by known patterns of binding.

2.4 Distribution of Signal: Distribution of coverage depth across the genome

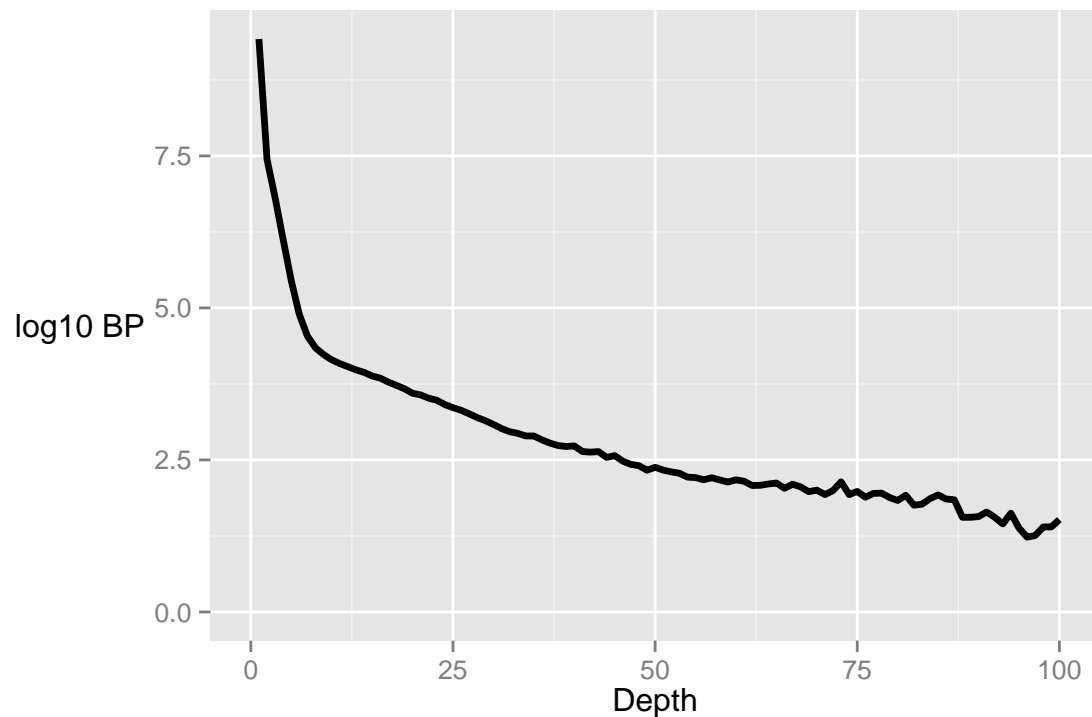
The final metrics to review are those of the distribution of global pile-up across the genome. We can access this in two ways within ChIPQC, first by visualising the histogram of coverage depths and secondly by applying the SSD metrics before and after removal of blacklisted regions.

First we draw the coverage histogram using `plotCoverageHistogram`. Note the cut-off at 100bp for visualisation purposes. If you want to replot the whole histogram you can extract the data using `coveragehistogram` function or extend x-axis as shown in Advanced Topics section.

```
coveragehistogram(exampleExp) [1:10]
##          0          1          2          3          4          5          6          7
## 2.644e+09 2.778e+07 6.486e+06 1.329e+06 2.805e+05 7.779e+04 3.432e+04 2.198e+04 1.709e+04
##          9
## 1.400e+04

sum(coveragehistogram(exampleExp) [-c(1:19)])
## [1] 43192

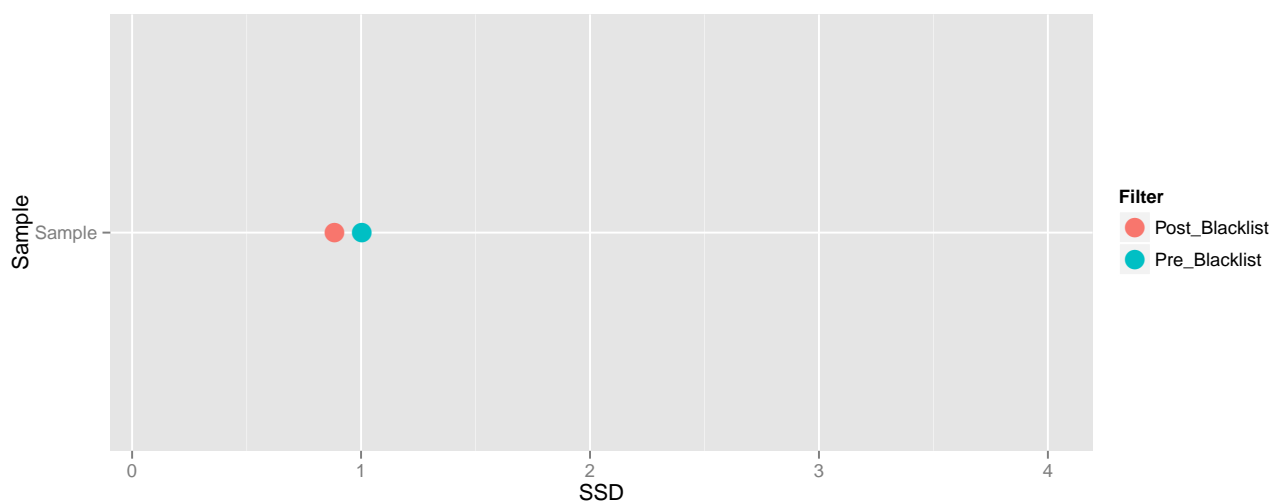
plotCoverageHist(exampleExp)
```



The coverage histogram shows that there is a significant stretch of high signal pile-up (more than 40,000 bps at greater than 20 depth). This may indicate significant signal associated to binding events but could also be from signal seen within blacklisted regions.

To assess the contribution of artefact signal to global distribution of signal pileup we measure SSD before and after exclusion of signal from known blacklisted regions.

```
plotSSD(exampleExp)
```



Here we find that the SSD signal is not greatly affected by blacklisted regions and so, taken together with the coverage histogram, indicates the Ebf1 ChIP has a clear ChIP-signal above background.

2.5 Conclusion

From the combination of metrics used here we can establish that this Ebf1 ChIP showed an enrichment for structured ChIP-signal, had a signal pile-up above that seen within artefact regions and that this ChIP is positively associated with TSS regions.

3 Assessing a ChIP-seq experiment using ChIPQC()

When assessing ChIP-seq quality it is most useful to consider your sample quality alongside other ChIP and input samples. Taken together, a full experiment of ChIP and input samples allows for the identification of expected enrichment of sample metrics above input but also the variation of sample quality between biological replicates and the identification of bias within input/control samples.

The ChIPQC function wraps the functionality of ChIPQCsample to allow for the assessment of within experiment ChIP sample/input quality alongside user supplied experimental metadata.

3.1 An example ChIP-seq experiment using ChIPQC()

The ChIPQC function accepts a samplesheet of metadata and file locations alongside the same set of arguments as ChIPQCsample.

An example of the layout for a sample sheet can be seen in the SampleSheet object below and in the ChIPQC and Diffbind vignettes. ChIPQC can also accept the DBA object from Diffbind package as a starting point for quality control. The result of a call to ChIPQC is the ChIPQCexperiment object which contains the list of ChIPQCsample objects.

```
load("/data/ChIPQC/SampleSheet.RData")
SampleSheet[1:3,]
#resExperiment = ChIPQC(SampleSheet,peaks=peaksFile,annotation="mm9",blacklist=BlackList)
load("/data/ChIPQC/BCell_Examples.RData")
resExperiment
```

```
##      SampleID Tissue  Factor Condition Treatment Peak caller      ControlID Replicate
## 1      DNase   Ch12   DNase          macs      Input_Ch12         1
## 2      Ebf1   ProB   Ebf1          macs Input_2_proB         1
## 3 H3K4me3_1  ProB H3K4me3          macs Input_2_proB         1
##                                     bamRead                                     bamControl
## 1 //AlignedData/DNaseDupMarked.bam //AlignedData/Input_Ch12DupMarked.bam
## 2 //AlignedData/Ebf1DupMarked.bam //AlignedData/Input_2_proBDupMarked.bam
## 3 //AlignedData/H3K4me3_1DupMarked.bam //AlignedData/Input_2_proBDupMarked.bam
##                                     Peaks
## 1 //AlignedData/DNaseDupMarked.bed
## 2 //AlignedData/Ebf1DupMarked.bed
## 3 //AlignedData/H3K4me3_1DupMarked.bed
```

##	Reads	Map%	Filt%	Dup%	ReadL	FragL	RelCC	SSD	RiP%	RiBL%
## H3K9ac_IkNeg	64585441	100	14.8	13.0	36	215	3.930	2.28	36.30	6.73
## H3K9ac_IkPos	69136473	100	13.9	14.7	36	213	4.170	2.01	37.90	5.84
## Ikaros_1_DPT	48253184	100	23.0	15.3	36	173	1.800	3.19	5.31	10.50
## Ikaros_1_preproB	35903381	100	24.7	15.3	36	175	0.278	3.43	3.50	12.50
## Ikaros_1_proB	24827255	100	23.2	13.1	36	165	0.921	2.28	3.49	10.50
## Ikaros_2_preproB	35493401	100	24.7	15.2	36	181	0.311	3.42	3.43	12.60

The QCmetric function now displays a table of metrics as seen for ChIPQCsample and similarly all accessors and plotting functions used for ChIPQCsample objects can be used with the ChIPQCexperiment object.

In addition to standard plotting routine, ChIPQCexperiment plots can be grouped by the metadata provided using the argument facetBy and for plotCoverageHistogram and plotCC methods the colours and line types controlled by colourBy and lineBy respectively.

To group/colour/line type by metadata, a character vector of the metadata column title/s to use may be provided.

```
facetBy = c("Tissue", "Factor", "Condition")
colourBy = c("Treatment", "Replicate")
lineBy = c("Replicate")
```

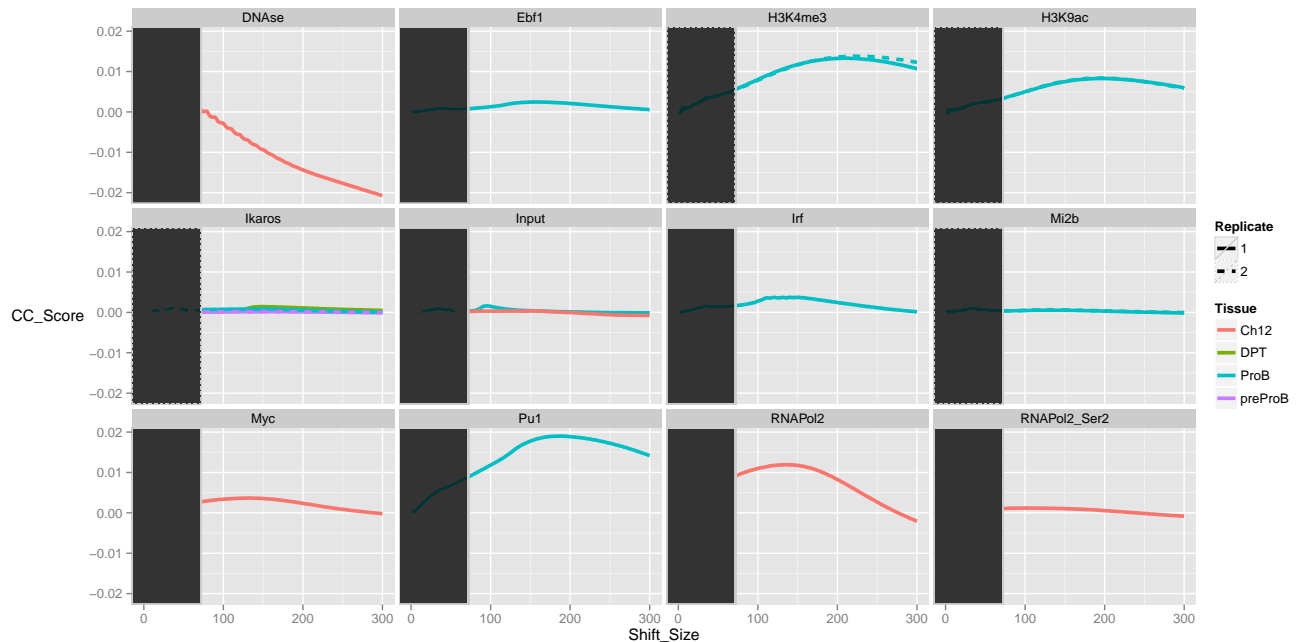
3.2 Examining Cross-coverage and FragCC/RelCC scores across an experiment

As with the ChIPQCsample object, the cross coverage scores for a group of samples can be plotted using plotCC(). By default samples will be groups by their "Tissue" and "Factor" combinations and coloured by their "Replicate" number. Here we group by "Factor", colour by "Tissue" and set the line type by the "Replicate" number.

```
FragmentLengthCrossCoverage(resExperiment)[1:3]
##          DNase          Ebf1 H3K4me3_IkNeg
## 0.0001494 0.0024607 0.0133150

RelativeCrossCoverage(resExperiment)[1:3]
##          DNase          Ebf1 H3K4me3_IkNeg
## 0.9277      2.6369      3.9340

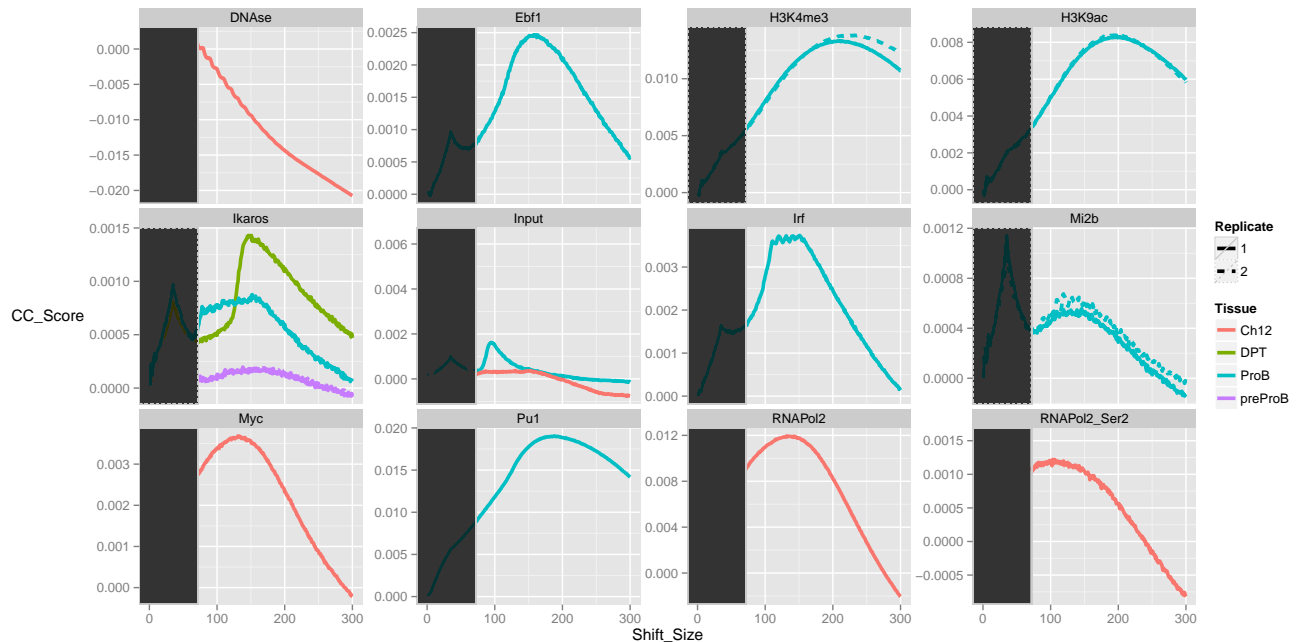
plotCC(resExperiment, colourBy="Tissue", facetBy="Factor", lineBy="Replicate")
```



From this, it is immediately apparent that some samples not only have much higher scores, and hence efficiency, than others but that their fragment lengths appear to be very different from each other. The DNase sample for example has a fragment length less than half of the Pu1 sample.

Now we have established the difference in total efficiency, we can look at the overall shape of cross-coverage scores and the relationship between signal peaks and artefact peaks in the cross-coverage scores. To help better visualise, we will apply a further facet wrap to the ggplot2 object returned by `plotCC` in order to compare within factors.

```
ccplot <- plotCC(resExperiment, colourBy="Tissue", facetBy="Factor", lineBy="Replicate")
ccplot + facet_wrap(~Factor, scales="free_y")
```



The free scaled cross-coverage score plots now reveal more about the distribution of signal within the samples. The Ebf1, Ikaros, Myc, Irf and RNA Pol2 all show tight peaks within their cross-coverage score profiles illustrating their sharp binding profiles as expected for a transcription factor (and RNA Pol2 around TSS/Enhancers.) The histone marks and Pu1 however show longer more diffused peaks reflecting the wider breadth of signal seen for these epigenetic marks.

The signal of the Ikaros ChIP between cell lines can also be seen to be highly variable with DP thymocytes containing highest RelCC scores, ProB lower and preProB the lowest. This reflects the increased concentrations of Ikaros along haemopoetic differentiation with DP thymocytes having the highest Ikaros levels.

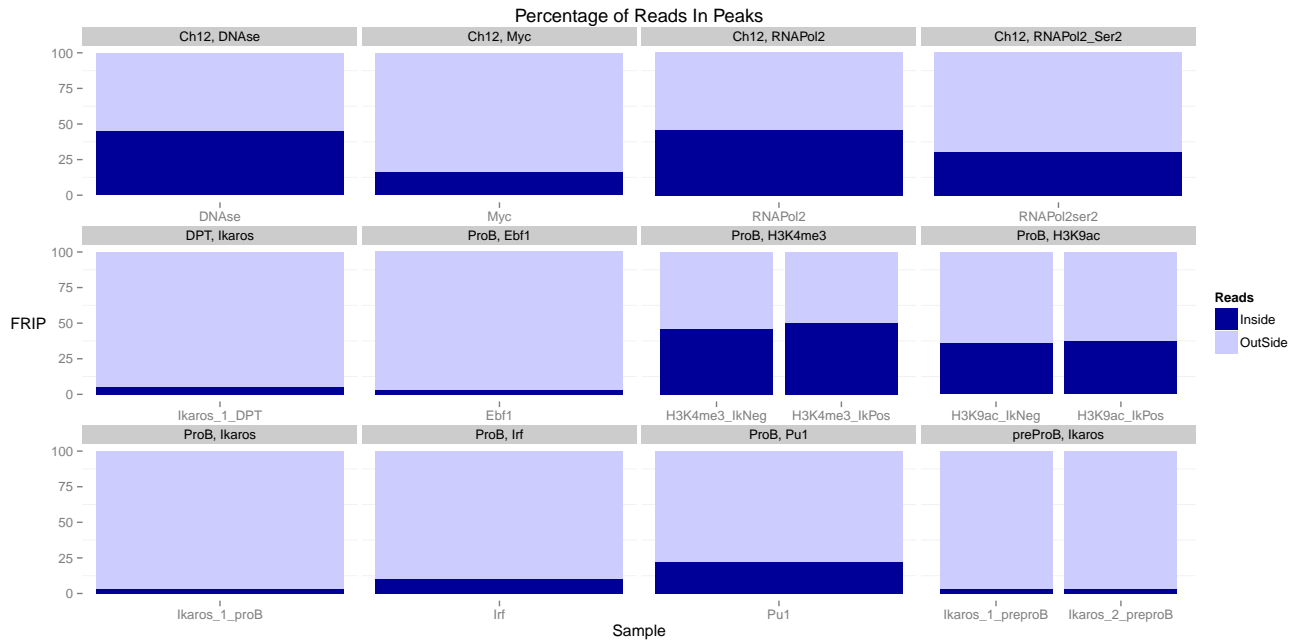
Finally, an enrichment for fragment length signal can be seen in the ProB input. The sharpness of this enrichment suggests a highly duplicated peak like signal within this track which would be cause for further investigation. This may be from the "sono-seq" effect commonly observed in input, where gel-selection of fragment lengths for input causes a small peak in cross coverage scores close to selected fragment length.

3.3 Distribution of Signal across a ChIPQC experiment

As with the the ChIPQCsample object, the fraction of signal in peaks, blacklists and annotated genomic intervals can provide an understanding of the ChIPs' efficiency and pattern of enrichment.

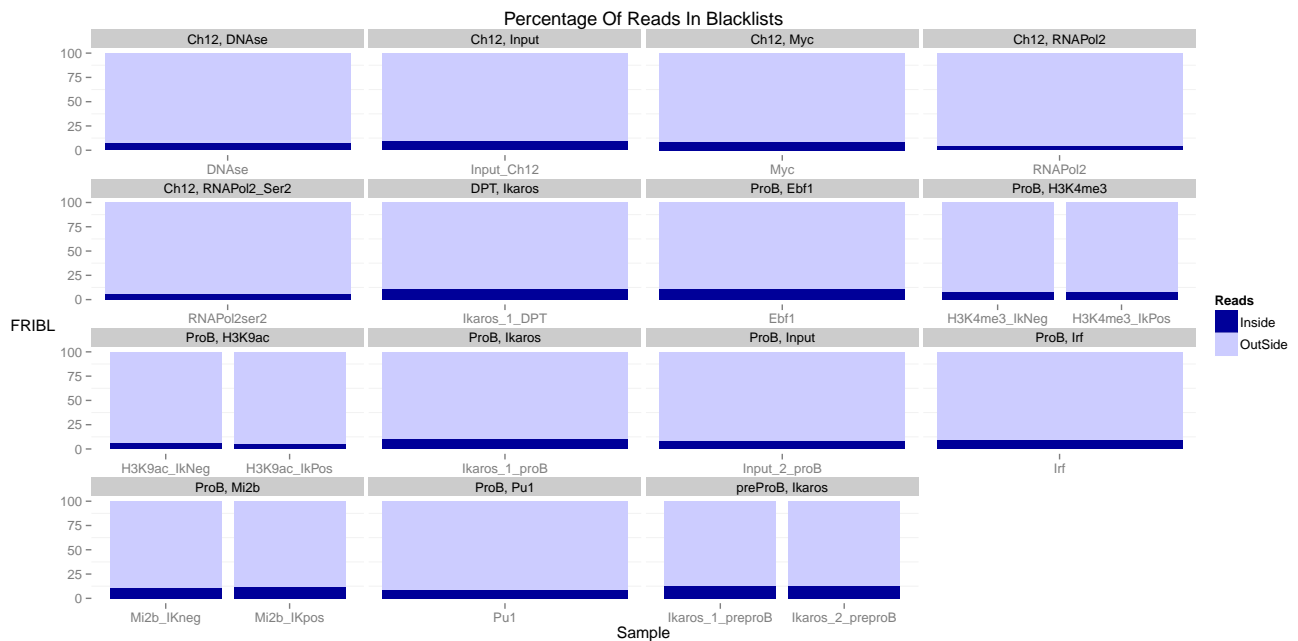
```
plotFrip(resExperiment)
```

```
## Using Sample as id variables
```



```
plotFribl(resExperiment)
```

```
## Using Sample as id variables
```



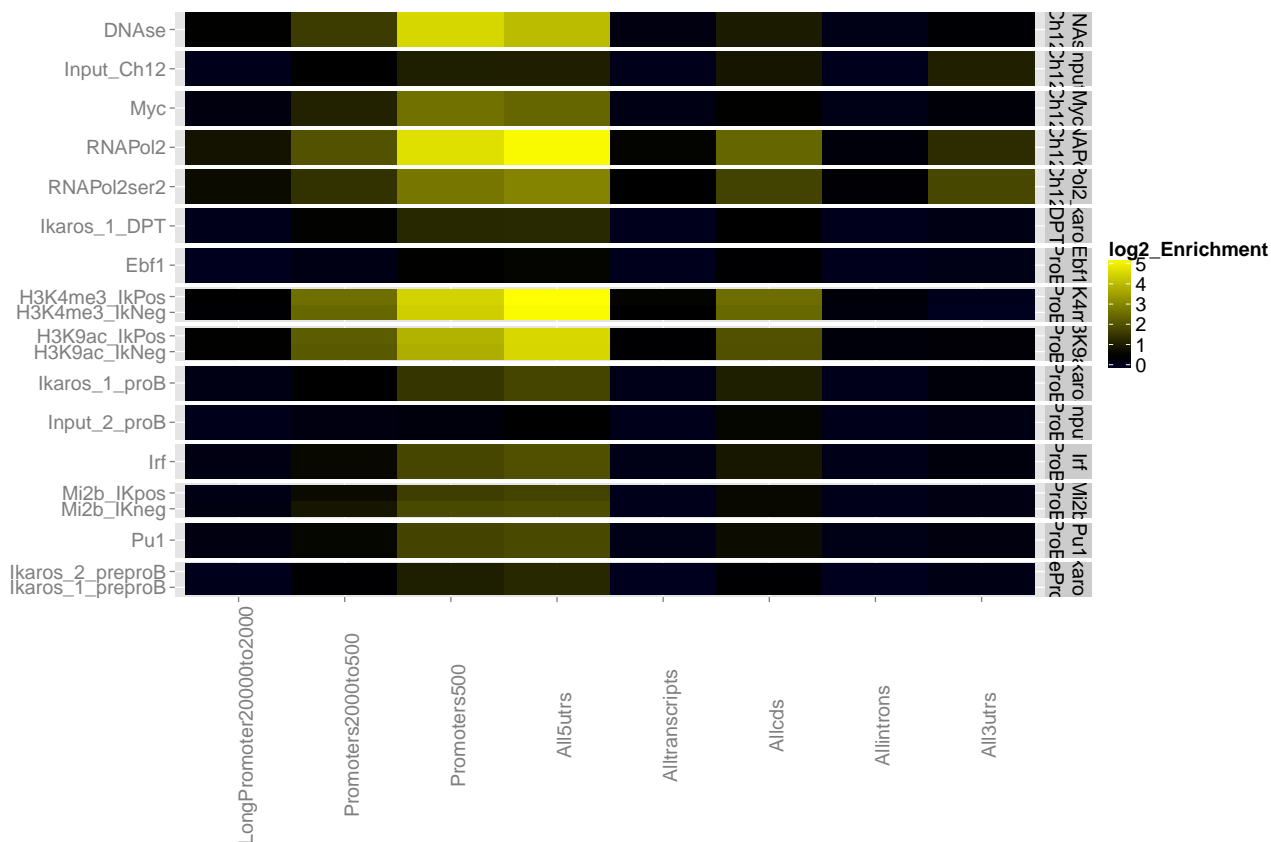
The output from `plotFrip` immediately identifies the Histone, RNA Pol2 and DNase ChIP as having the highest enrichment (25% to 50%) for reads in peaks as expected for such epigenetic marks. Also apparent are the significant enrichment seen within Pu1, Myc and Irf whereas the Ikaros ChIPs all show considerably lower enrichment for signal within peaks.

The Fribl plot here shows that all samples have equivalent levels of signal so no outlier or need to

investigate signal within known blacklisted regions.

Samples such as RNA Pol2, DNase and Histone marks will have an expected enrichment for genomic locations. Here RNA Pol2 should be expected to have a stonger enrichment at the TSS than RNA Pol2ser2 where as RNA Pol2ser2 should show enrichment within the 3'UTRs regions.

```
plotRegi(resExperiment)
```



```
regi(resExperiment) ["All3utrs",]
```

##	DNase	Ebf1	H3K4me3_IkNeg	H3K4me3_IkPos	H3K9ac_IkNeg
##	0.42145	0.08548	-0.08227	0.01885	0.39312
##	H3K9ac_IkPos	Ikaros_1_DPT	Ikaros_1_preproB	Ikaros_1_proB	Ikaros_2_preproB
##	0.40440	0.10945	0.10805	0.29672	0.10833
##	Input_2_proB	Input_Ch12	Irf	Mi2b_IKneg	Mi2b_IKpos
##	0.15272	1.08960	0.28397	0.16221	0.16271
##	Myc	Pu1	RNAPol2	RNAPol2ser2	
##	0.32845	0.23294	1.31672	1.80521	

From the Regi plot it can be seen that all histone marks, DNA and RNA Pol2 show the expected enrichment across gene regions. Combined with the output from the reg function, the RNA Pol2Ser2 has the greatest enrichment at 3'UTRs and so the expected pattern of enrichment.

To better visualize this enrichment we can adjust the scale to the enrichment seen in Ch12 input for

3'UTRs.

```
plotRegi(resExperiment)+scale_fill_gradient2(low="white",high="red",
  mid="white",midpoint=regi(resExperiment)["All3utrs","Input_Ch12"])
```

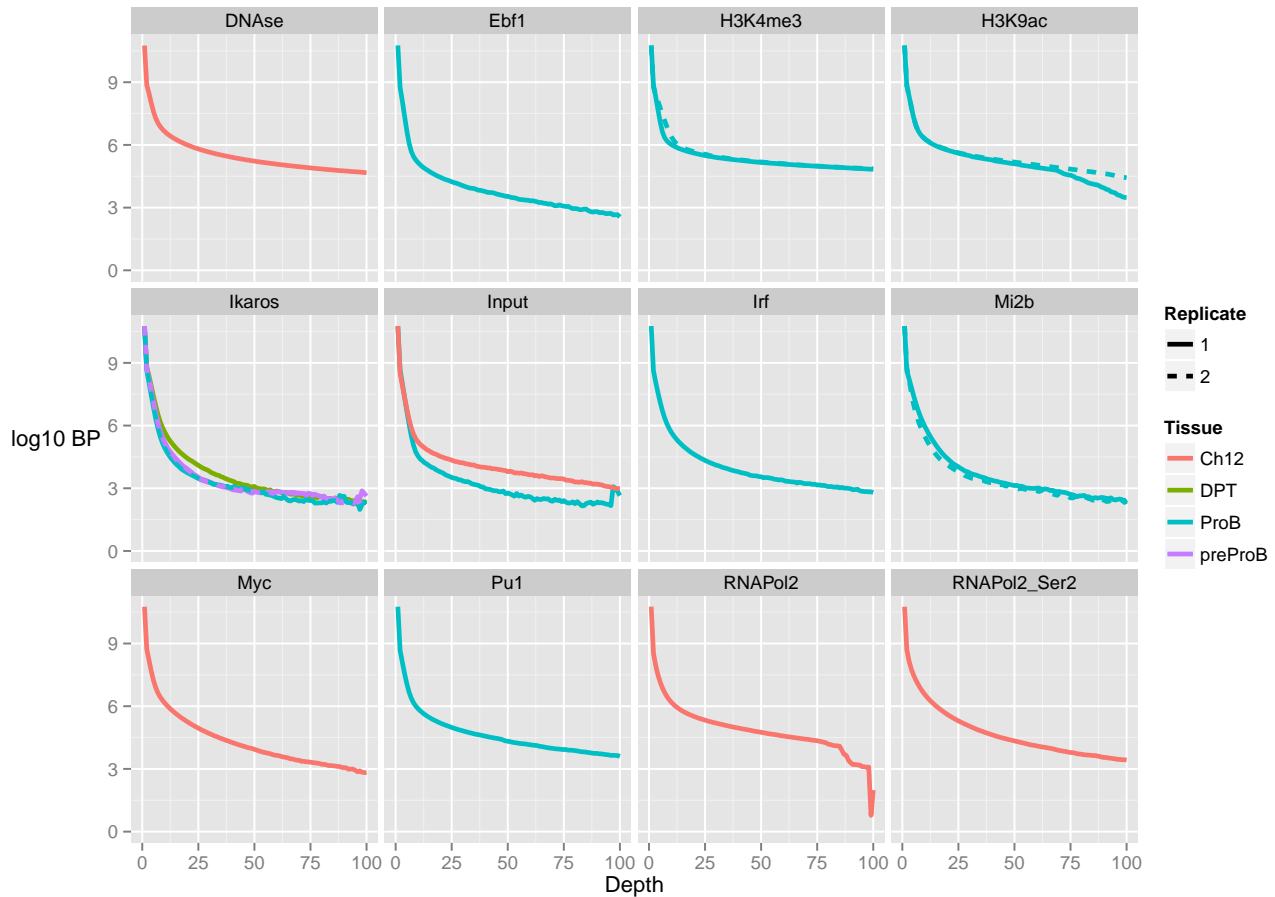
Scale for 'fill' is already present. Adding another scale for 'fill', which will replace the existing scale.



Both inputs showed a small but comparatively low enrichment for reads in genic regions. Such enrichment around gene regions can typically be seen for input samples due to the increased accessibility of chromatin to fragmentation around TSSs.

As with the *ChIPQC*sample object we can plot the coverage histogram and SSD before and after blacklisting by using the *plotCoverageHist* and *plotSSD* functions.

```
plotCoverageHist(resExperiment, facetBy="Factor", colourBy="Tissue", lineBy="Replicate")
```

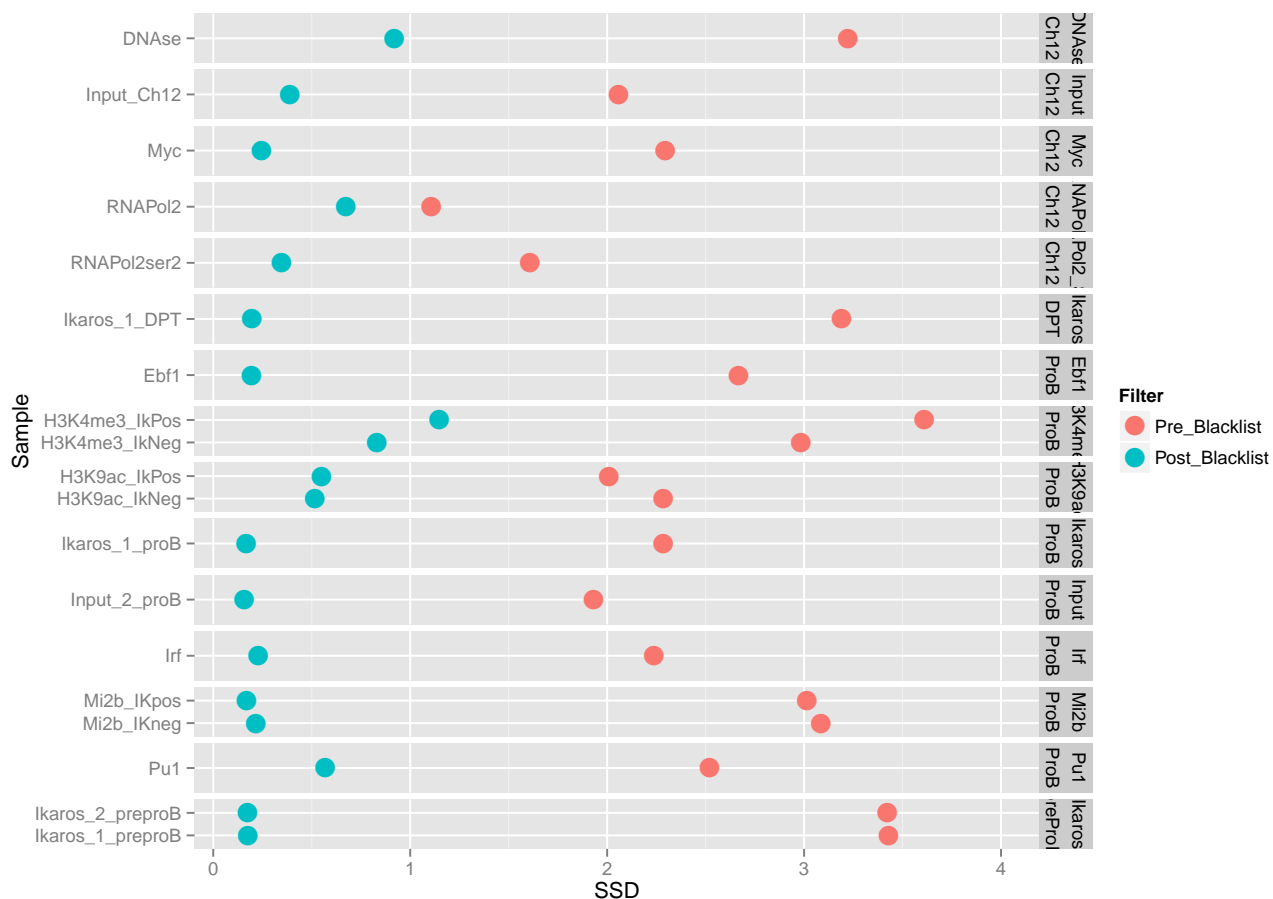


The coverage histogram shows the expected greater spans of high signal in Histone, RNA Pol2 and DNase ChIPs as well as for transcription factors with both high Rip% and RelCC scores. Although most ProB transcription factors have greater spans of high depth than their input, enrichment for these transcription factors can be seen to be much smaller than observed for Histone, RNA Pol2 and DNase ChIPs.

The Ikaros ChIPs universally had low enrichment and as seen with other metrics DP thymocytes had the greatest signal.

The two inputs samples show very different patterns of signal depth. The Ch12 input show considerable span of high signal where as the ProB shows a spike in signal at above 98 reads high, in keeping with observed high duplication rate for this sample and fragment length peak in cross-coverage scores.

```
plotSSD(resExperiment)
## Using Sample as id variables
```

The plot of SSD scores before and after blacklisting show that the effect of blacklisting on SSD scores is dramatic.

As with the coverage histogram plots, the histone, polymerase, DNase marks as well as high scoring TFs have greatest SSDs after blacklisting.

The SSD for Ikaros after blacklisting can be seen to be just above background with again DP thymocytes having the greatest score.

For the inputs, the ProB input can be seen to have its SSD score reduced to a background level of around 0.14 indicating the successful removal of artefact signal. The Ch12 input however can be seen to not drop to the background level after blacklisting suggesting remaining regions of artefact signal. The failure to reduce SSD after blacklisting suggests the persistent presence of artefacts within the Ch12 input and flags this sample for further blacklisting.

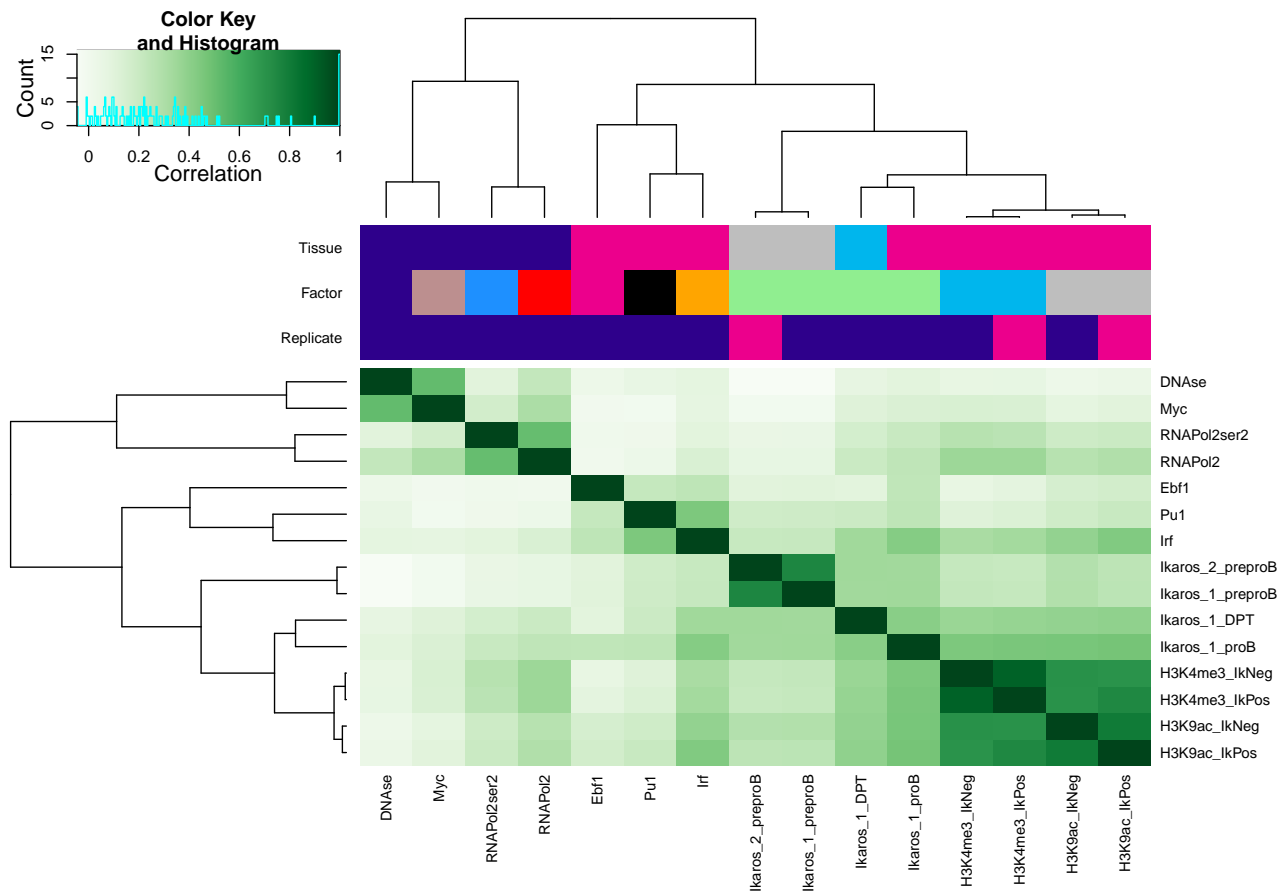
3.4 Assessing sample similarity with Diffbind

A final set of metrics useful for ChIP-quality relate to the correlation between binding events across samples within a ChIPQC experiment.

The Diffbind package allows for the clustering of samples based on the co-occurrence of peaks within

samples. When analysing an experiment ChIPQC will perform a sample clustering by default as well as construct the Diffbind DBA object. To produce the a sample heatmap, the `plotCorHeatmap` function can be used.

```
plotCorHeatmap(resExperiment, facetBy="Factor")
```



The clustering and heatmaps generated from `plotCorHeatmap` allow us to identify which samples are most closely related and so both the reproducibility of replicates and the similarity in binding profiles of different epigenetic marks.

Here we see that all replicates cluster tightly together, illustrating the relative reproducibility of our replicates. Further to this we see that the samples broadly group into their respective tissues with DNase, RNA Pol2 and Myc forming the Ch12 cluster.

The ProB transcription factor (Pu1, Irf and Ebf) are seen to group together as expected and are found to cluster away from ProB histone marks. This suggests that these transcription factors may be less associated to these histone marks than the Ikaros ChIPs.

3.5 Conclusion

The analysis of the quality of this experiment's ChIP quality using ChIPQC has identified several informatics characteristics of the data as well as highlighted variability in quality across ChIPs and cell-lines.

The Histone , RNA Pol2 and DNase ChIPs all have high RIP% and SSD after blacklisting as expected from broader epigenetic marks. Cross-coverage score profiles illustrate the broader regions of enrichment for the Histone marks, tighter profiles for the RNA Pol2 and RNA Pol2-ser2 marks, due to their narrow enrichment in TSSs and a sharp narrow profile for The DNase ChIP. These epigenetic marks showed the characteristic enrichment within TSSs with RNAPol2-ser2 most enriched for 3'UTRs.

The transcription factor ChIPs showed a much wider variability in RIP%, SSD and RelCC than seen for histone,pol2 and DNase ChIPs. Pu1 and Irf were found to be highly efficient ChIPs and Pu1 was seen in its cross coverage scores to have a broader enrichment pattern than seen for other TFs. The Ikaros ChIPs were found to have acceptable but low enrichment for signal by all metrics, and the enrichment for Ikaros ChIP signal was seen to fit known concentrations of Ikaros within these cell-lines. The ProB Pu1, Ebf and Irf transcription factors were found to cluster away from Ikaros ChIPs indicating a greater co-occurrence of binding among them than with Ikaros ChIPs.

The inputs used in this study showed different sources of artefact contamination. The Ch12 input showed a strong artefact, read length, peak in cross-coverage scores and significant pile-up of signal. Following blacklisting it's SSD score didn't drop to that of a background level and so flagged this as a control for further blacklisting. The ProB input however showed peak like signal in its cross-coverage scores profile, a high level of duplication and a spike in its coverage histogram but removal of known blacklisted regions removed much of artefact signal as measured by SSD. Taken together, this suggests that the ProB input contains highly duplicated peak shaped spikes in signal which were contained within blacklisted regions.

4 Advanced Topics

4.1 Providing additional data to ChIPQC plotting and reporting

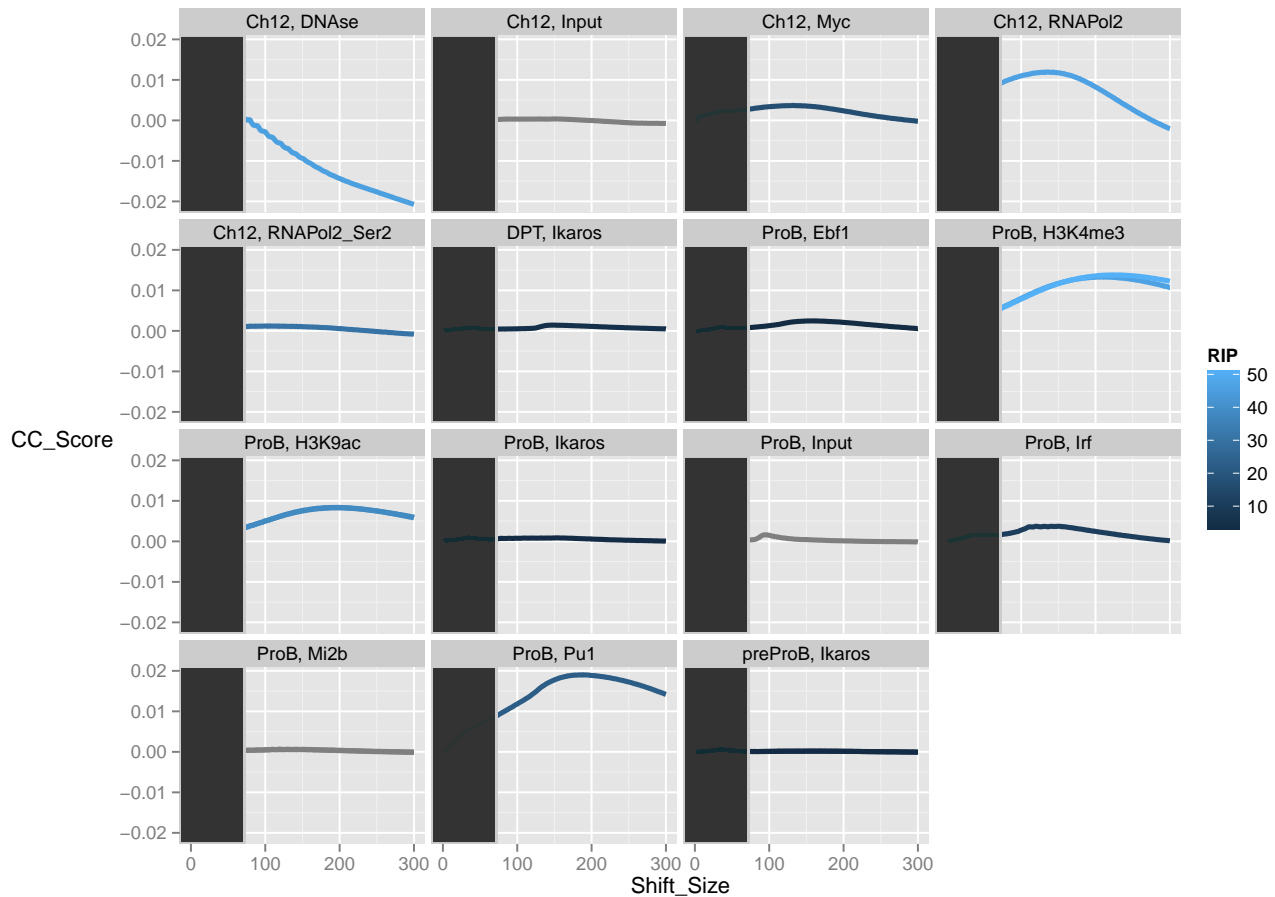
The sample sheets for ChIPQC may contain the optional metadata columns "Tissue", "Factor", "Condition" and "Treatment". In order to allow the user to specify custom metadata for their plotting and reporting, the additional `addMetadata` argument can be supplied with a data frame of sample IDs and associated additional metadata. The first column for `addMetadata` data frame must be `SampleID` and remaining columns may be categorical or discrete data.

Here, first we illustrate the relationship between RIP and cross coverage scores by including RIP metrics as metadata and setting the column as `colourBy` argument.

Then we can use `addMetadata` argument to group SSD by their ChIP type to highlight higher SSD scores typically seen to Epigenetic marks.

```
metrics <- QCmetrics(resExperiment)
metricsMetadata <- data.frame(SampleID=rownames(metrics),RIP =metrics[, "RiP%",drop=T])
```

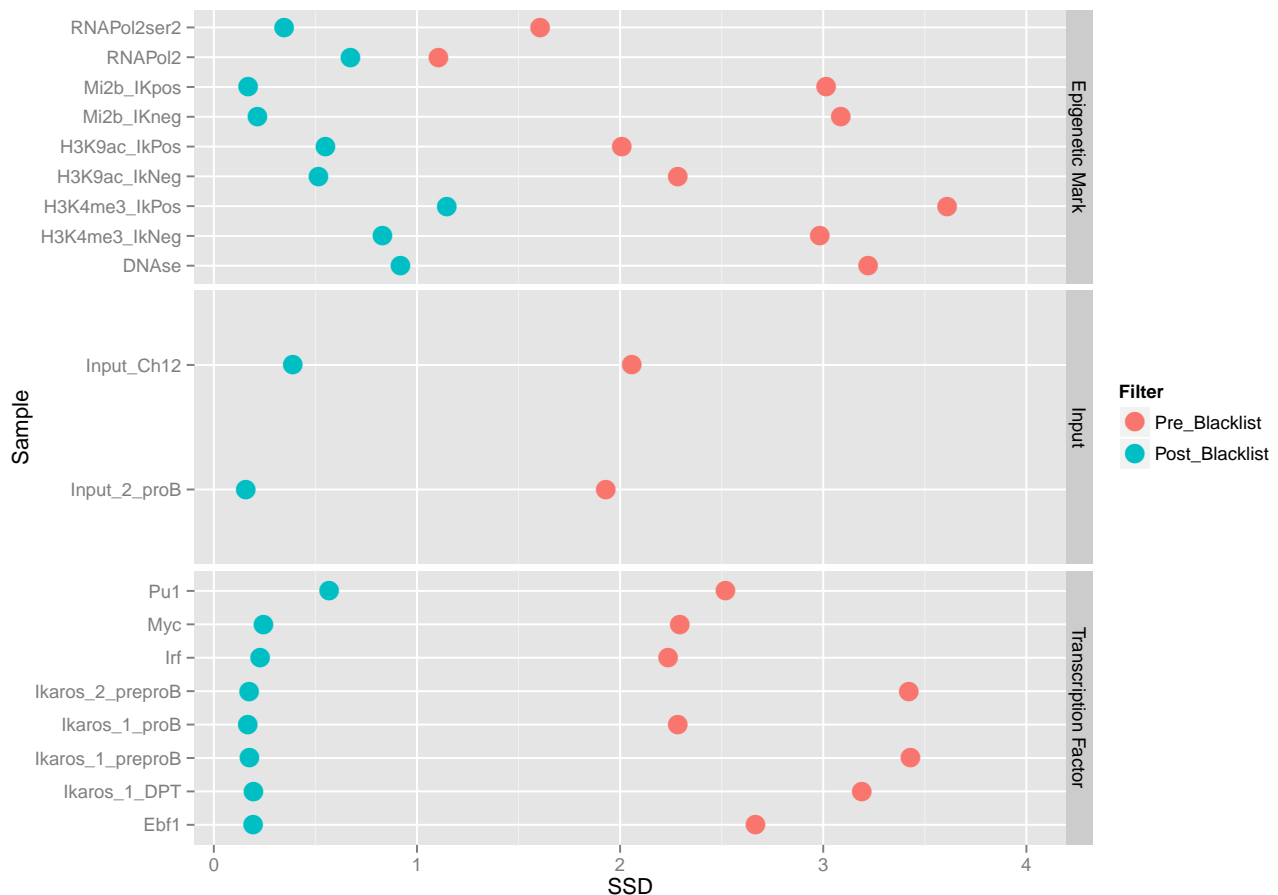
```
plotCC(resExperiment, addMetaData = metricsMetadata, colourBy="RIP")
```



```
metricsMetadata <- data.frame(SampleID=rownames(metrics),ChIPType=
  c(rep("Epigenetic Mark",1),
    rep("Transcription Factor",1),
    rep("Epigenetic Mark",4),
    rep("Transcription Factor",4),
    rep("Input",2),
    rep("Transcription Factor",1),
    rep("Epigenetic Mark",2),
    rep("Transcription Factor",2),
    rep("Epigenetic Mark",2)))
```

```
plotSSD(resExperiment, addMetaData = metricsMetadata, facetBy="ChIPType")
```

```
## Using Sample as id variables
```



4.2 Specifying custom annotation

ChIPQC produces a simple metric for enrichment of signal in known genomic annotation. By default annotation is provided from the TranscriptDB packages specified by "species" argument. ChIPQC also allows for the specification of genomic annotation in the form of GRanges objects and so enrichment in user defined regions can be assessed.

Annotation must be provided as a named list with the first element being "version" and the remaining list elements being GRanges objects.

Here we will run ChIPQCsample function with DNase and Histone mark peaks as the annotation.

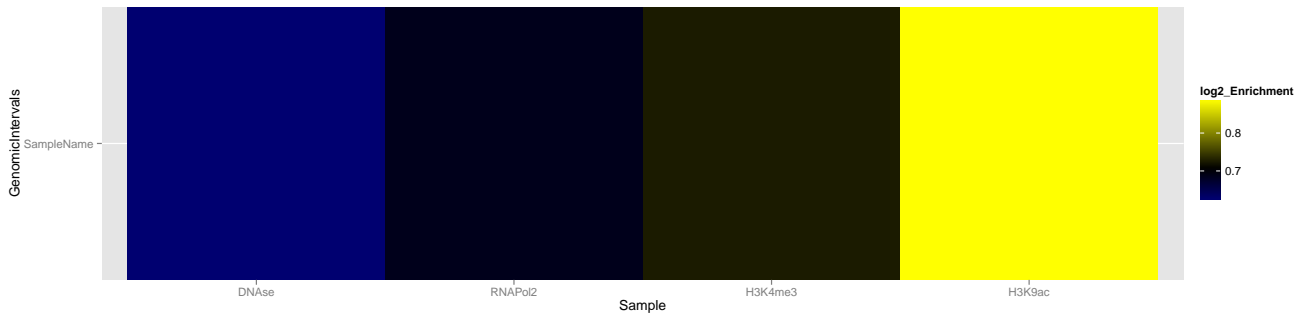
```
DNasePeaks <- peaks(QCsample(resExperiment, "DNase"))
H3K4me3Peaks <- peaks(QCsample(resExperiment, "H3K4me3_IkPos"))
H3K9acPeaks <- peaks(QCsample(resExperiment, "H3K9ac_IkPos"))
RNAPol2Peaks <- peaks(QCsample(resExperiment, "RNAPol2"))

customAnnotation <- list(version="custom",
                          DNase=DNasePeaks,
                          RNAPol2=RNAPol2Peaks,
```

```

H3K4me3=H3K4me3Peaks,
H3K9ac=H3K9acPeaks)
#bamFile <- "/data/ChIPQC/Chr11_Ebf1DupMarked.bam
bamFile <- "/data/ChIPQC/Chr11_Ebf1DupMarked.bam"
#exampleExpCA = ChIPQCsample(bamFile,peaks=NULL,annotation=customAnnotation,chromosomes=
load("/data/ChIPQC/exampleCustomAnnotation.RData")
plotRegi(exampleExpCA)

```



The `plotRegi` now shows the enrichment over expected for the custom annotation. The order of display of custom annotation is dictated by the order in annotation list and so may be rearranged as user desires.

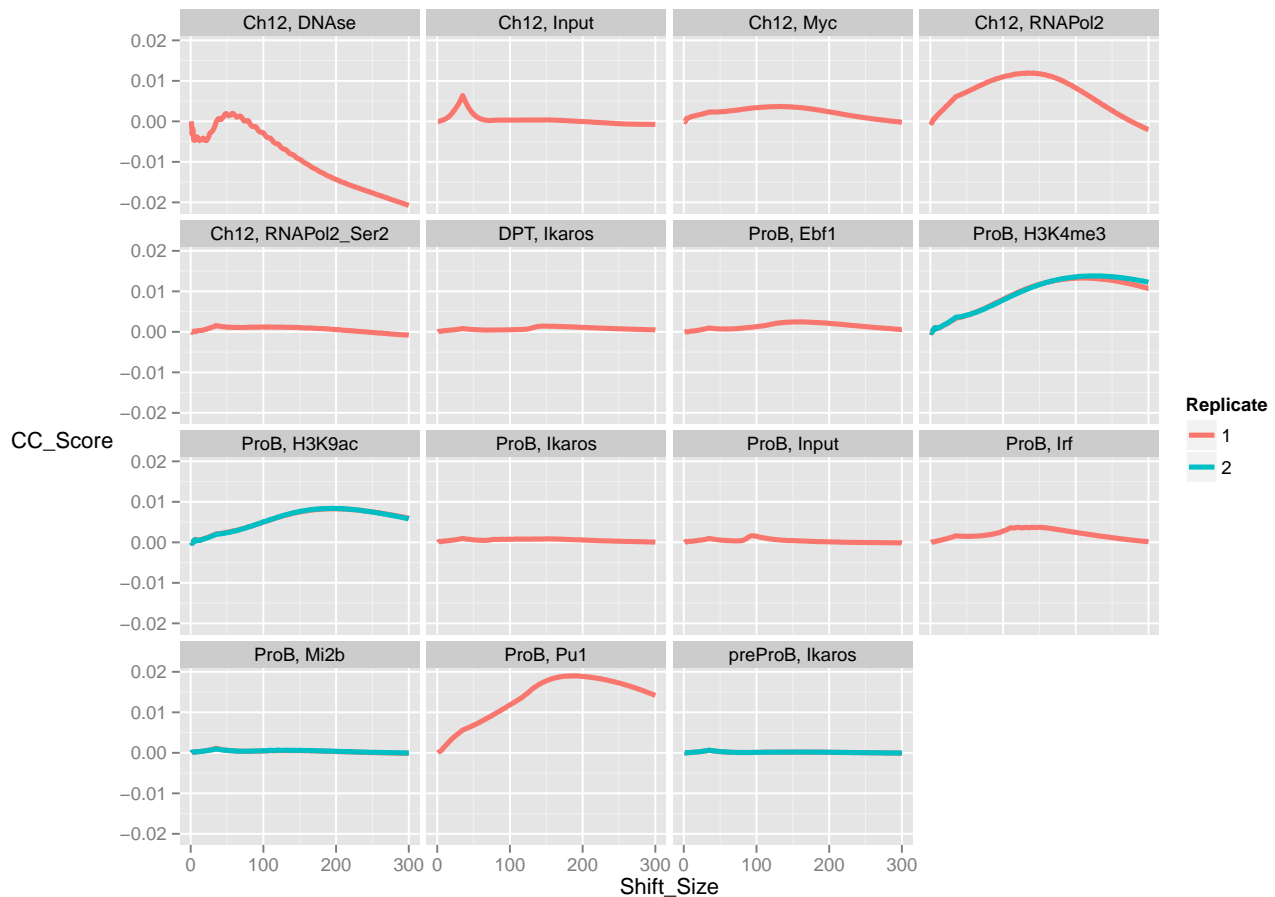
4.3 Some plotting tips

The `plotCC` function produces a line plot with the section excluded from identification of fragment length cross-coverage score shaded in grey. To remove this the `plotCC` ggplot object maybe altered and replotted.

```

ccplot <- plotCC(resExperiment)
ccplot$layers <- ccplot$layers[1]
ccplot

```



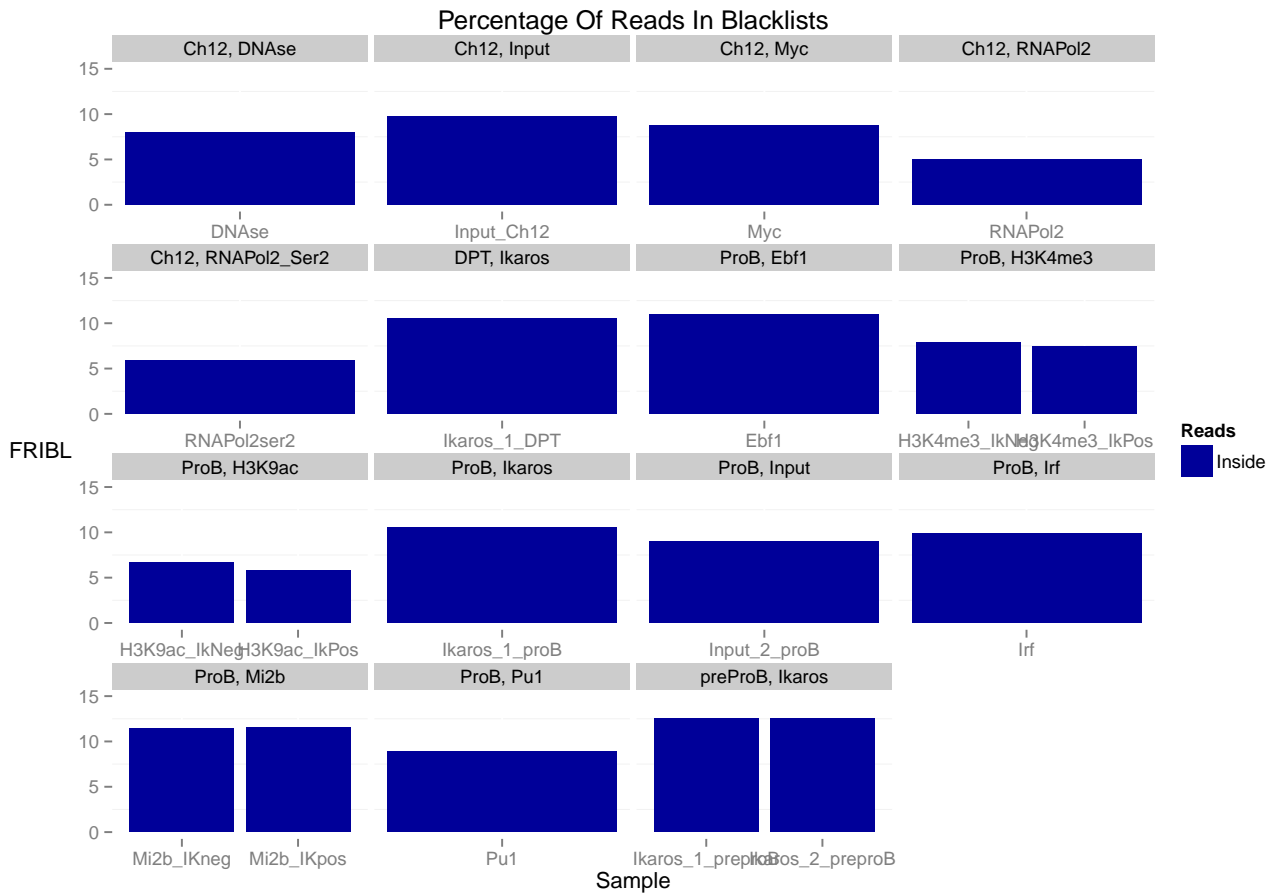
Another useful tip is expanding or narrowing the Y and X axis limits post plotting. This can be done by simply adding the `xlim()` and `ylim()` arguments

```
plotFribl(resExperiment)+ylim(0,15)
```

```
## Using Sample as id variables
```

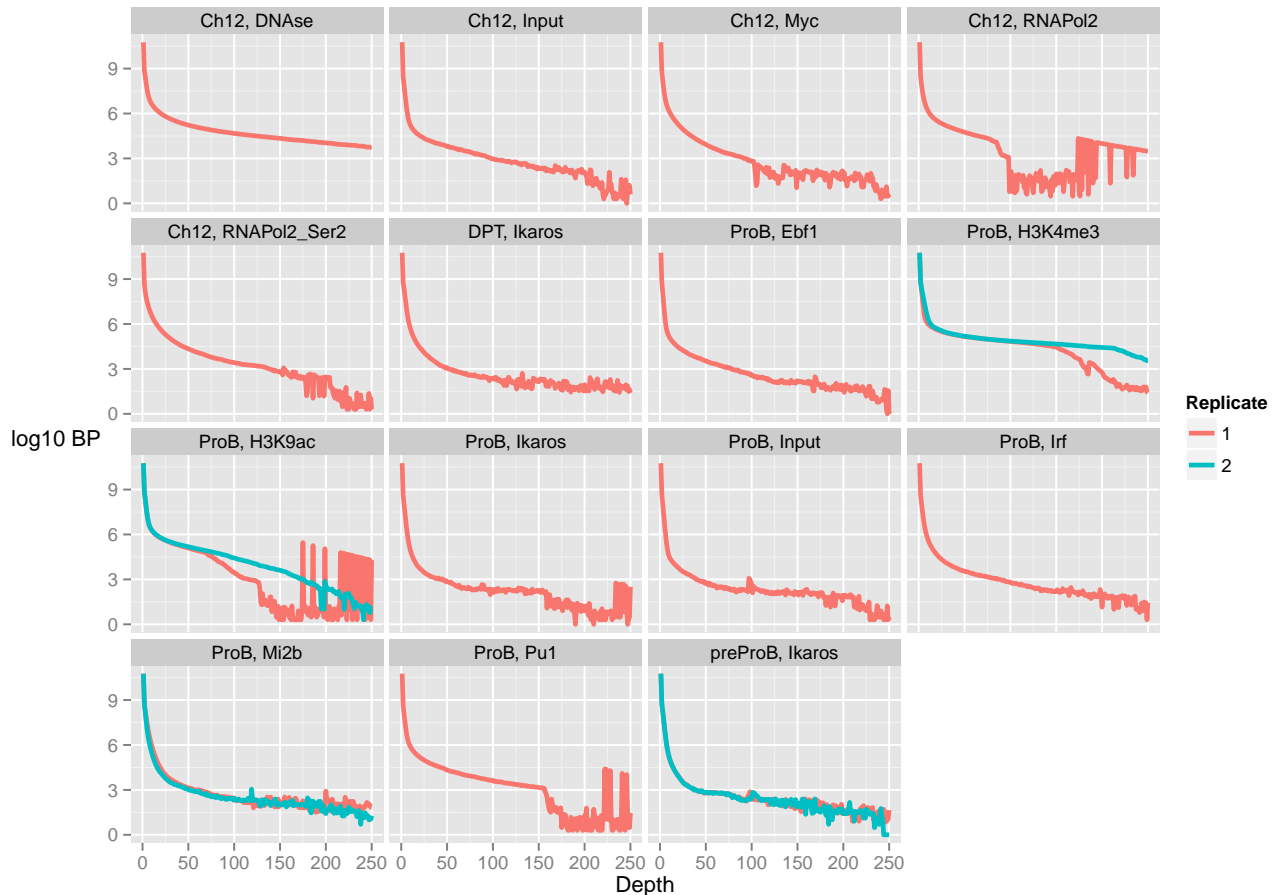
```
## Warning: Removed 1 rows containing missing values (position_stack).
## Warning: Removed 1 rows containing missing values (position_stack).
## Warning: Removed 1 rows containing missing values (position_stack).
## Warning: Removed 1 rows containing missing values (position_stack).
## Warning: Removed 1 rows containing missing values (position_stack).
## Warning: Removed 1 rows containing missing values (position_stack).
## Warning: Removed 1 rows containing missing values (position_stack).
## Warning: Removed 2 rows containing missing values (position_stack).
## Warning: Removed 2 rows containing missing values (position_stack).
## Warning: Removed 1 rows containing missing values (position_stack).
## Warning: Removed 1 rows containing missing values (position_stack).
## Warning: Removed 1 rows containing missing values (position_stack).
## Warning: Removed 2 rows containing missing values (position_stack).
## Warning: Removed 1 rows containing missing values (position_stack).
```

```
## Warning: Removed 2 rows containing missing values (position_stack).
```



```
plotCoverageHist(resExperiment)+xlim(0,250)
```

Scale for 'x' is already present. Adding another scale for 'x', which will replace the existing scale.



5 Analysing the Tamoxifen dataset

Now, given what we have seen with the previous experiment, let's take a look at the tamoxifen dataset.

```
data(tamoxifen_QC)
QCmetrics(tamoxifen)
```

##	Reads	Map%	Filt%	Dup%	ReadL	FragL	RelCC	SSD	RiP%	RiBL%
## BT4741	776353	100	15.8	8.42	28	153	2.110	1.88	14.90	1.75
## BT4742	782419	100	15.1	10.30	28	147	2.010	1.63	13.80	1.68
## MCF71	438994	100	20.9	21.30	28	134	2.300	2.52	26.50	1.70
## MCF72	465700	100	20.8	4.84	28	155	2.090	1.65	16.10	2.17
## MCF73	577273	100	19.1	10.30	28	153	2.460	2.18	21.80	1.83
## T47D1	507492	100	21.1	7.86	28	153	1.520	1.59	9.62	2.24
## T47D2	1831766	100	19.3	9.56	28	162	1.530	2.30	5.70	2.26
## TAMR1	747610	100	17.4	15.50	28	151	2.490	2.15	19.30	2.05
## TAMR2	728601	100	18.4	5.84	28	149	1.850	1.45	12.10	1.56
## ZR751	804427	100	12.0	15.80	28	158	2.950	3.52	30.70	1.39

##	ZR752	2918549	100	11.6	23.80	28	160	3.140	4.45	20.70	1.22
##	BT474c	598010	100	18.1	3.30	28	105	0.202	1.14	2.98	1.72
##	MCF7c	485192	100	26.3	1.70	28	109	0.108	1.40	2.65	2.47
##	T47Dc	400396	100	44.2	31.60	36	196	0.268	6.02	1.18	8.56
##	TAMRc	779102	100	19.3	6.16	28	110	0.267	1.38	2.21	2.06
##	ZR75c	1023987	100	26.2	20.10	36	220	0.573	5.36	1.50	5.35

With this data, take some time to identify key points about the dataset using ChIPQC.

1. Which samples have the highest enrichment?
2. Are sample replicates reproducible in terms of quality and peak occupancy?
3. Is enrichment consistent across samples?
4. Is there evidence of artefacts in input samples?
5. Are sources of artefacts the same?

We will look at this at the end of practical.

References

- [1] Carroll T Liang Z , Salama R, Stark R and Santiago I *Impact of artifact removal on ChIP quality metrics in ChIP-seq and ChIP-exo data*, *Frontiers in Genetics*, April 2014.

6 Session Info

- R version 3.1.0 (2014-04-10), x86_64-unknown-linux-gnu
- Locale: LC_CTYPE=en_US.UTF-8, LC_NUMERIC=C, LC_TIME=C, LC_COLLATE=C, LC_MONETARY=C, LC_MESSAGES=C, LC_PAPER=C, LC_NAME=C, LC_ADDRESS=C, LC_TELEPHONE=C, LC_MEASUREMENT=C, LC_IDENTIFICATION=C
- Base packages: base, datasets, grDevices, graphics, grid, methods, parallel, stats, utils
- Other packages: BSgenome 1.32.0, BiocGenerics 0.10.0, BiocInstaller 1.14.2, Biostrings 2.32.1, ChIPQC 1.0.9, DiffBind 1.10.2, GenomInfoDb 1.0.2, GenomicAlignments 1.0.3, GenomicRanges 1.16.3, IRanges 1.22.9, Rsamtools 1.16.1, XVector 0.4.0, ggplot2 1.0.0, knitr 1.6, limma 3.20.8
- Loaded via a namespace (and not attached): BBmisc 1.7, BatchJobs 1.3, Biobase 2.24.0, BiocParallel 0.6.1, BiocStyle 1.2.0, DBI 0.2-7, KernSmooth 2.23-12, MASS 7.3-33, Nozzle.R1 1.1-1, RColorBrewer 1.0-5, RSQLite 0.11.4, Rcpp 0.11.2, ShortRead 1.22.0, amap 0.8-12, bitops 1.0-6, brew 1.0-6, caTools 1.17, checkmate 1.1, chipseq 1.14.0, codetools 0.2-8, colorspace 1.2-4, digest 0.6.4, edgeR 3.6.7, evaluate 0.5.5, fail 1.2, foreach 1.4.2, formatR 0.10, gdata 2.13.3, gplots 2.14.1, gtable 0.1.2, gtools 3.4.1, highr 0.3, hwriter 1.3, iterators 1.0.7, labeling 0.2, lattice 0.20-29, latticeExtra 0.6-26, munsell 0.4.2, plyr 1.8.1, proto 0.3-10, reshape2 1.4, scales 0.2.4, sendmailR 1.1-2, stats4 3.1.0, stringr 0.6.2, tools 3.1.0, zlibbioc 1.10.0