

# ChIP-Seq workflow template: Some Descriptive Title

Project ID: ChIPseq-PI\_Name\_Organism\_Jun2015

Project PI: First Last (first.last@inst.edu)

Author of Report: First Last (first.last@inst.edu)

May 1, 2016

## Contents

---

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Background and objectives . . . . .	2
1.2	Experimental design . . . . .	2
<b>2</b>	<b>Load workflow environment</b>	<b>2</b>
2.1	Load packages and sample data . . . . .	2
2.2	Experiment definition provided by <code>targets</code> file . . . . .	2
<b>3</b>	<b>Read preprocessing</b>	<b>3</b>
3.1	Read quality filtering and trimming . . . . .	3
3.2	FASTQ quality report . . . . .	3
<b>4</b>	<b>Alignments</b>	<b>3</b>
4.1	Read mapping with <code>Bowtie2</code> . . . . .	4
4.2	Read and alignment stats . . . . .	4
4.3	Create symbolic links for viewing BAM files in IGV . . . . .	4
<b>5</b>	<b>Peak calling with MACS2</b>	<b>4</b>
5.1	Merge BAM files of replicates prior to peak calling . . . . .	4
5.2	Peak calling without input/reference sample . . . . .	5
5.3	Peak calling with input/reference sample . . . . .	5
<b>6</b>	<b>Annotate peaks with genomic context</b>	<b>5</b>
6.1	Annotation with <code>ChIPpeakAnno</code> package . . . . .	5
6.2	Annotation with <code>ChIPseeker</code> package . . . . .	6
<b>7</b>	<b>Count reads overlapping peak regions</b>	<b>6</b>
<b>8</b>	<b>Differential binding analysis of peaks</b>	<b>6</b>
<b>9</b>	<b>GO term enrichment analysis</b>	<b>6</b>
<b>10</b>	<b>Motif analysis</b>	<b>7</b>
10.1	Parse DNA sequences of peak regions from genome . . . . .	7
10.2	Motif discovery with <code>BCRANK</code> . . . . .	7
<b>11</b>	<b>Version Information</b>	<b>8</b>
<b>12</b>	<b>Funding</b>	<b>9</b>

## 13 References

9

## 1 Introduction

---

### 1.1 Background and objectives

This report describes the analysis of several ChIP-Seq experiments studying the DNA binding patterns of the transcriptions factors ... from *organism* ....

### 1.2 Experimental design

Typically, users want to specify here all information relevant for the analysis of their NGS study. This includes detailed descriptions of FASTQ files, experimental design, reference genome, gene annotations, etc.

## 2 Load workflow environment

---

### 2.1 Load packages and sample data

The *systemPipeR* package needs to be loaded to perform the analysis steps shown in this report ([Girke, 2014](#)).

```
library(systemPipeR)
```

Load workflow environment with sample data into your current working directory. The sample data are described [here](#).

```
library(systemPipeRdata)
genWorkenvir(workflow="chipseq")
setwd("chipseq")
```

In the workflow environments generated by *genWorkenvir* all data inputs are stored in a *data/* directory and all analysis results will be written to a separate *results/* directory, while the *systemPipeChIPseq.Rnw* script and the *targets* file are expected to be located in the parent directory. The R session is expected to run from this parent directory. Additional parameter files are stored under *param/*.

To work with real data, users want to organize their own data similarly and substitute all test data for their own data. To rerun an established workflow on new data, the initial *targets* file along with the corresponding FASTQ files are usually the only inputs the user needs to provide.

If applicable users can load custom functions not provided by *systemPipeR*. Skip this step if this is not the case.

```
source("systemPipeChIPseq_Fct.R")
```

### 2.2 Experiment definition provided by targets file

The *targets* file defines all FASTQ files and sample comparisons of the analysis workflow.

```
targetspath <- system.file("extdata", "targets_chip.txt", package="systemPipeR")
targets <- read.delim(targetspath, comment.char = "#")
targets[1:4,-c(5,6)]
```

	FileName	SampleName	Factor	SampleLong	SampleReference
1	./data/SRR446027_1.fastq	M1A	M1	Mock.1h.A	
2	./data/SRR446028_1.fastq	M1B	M1	Mock.1h.B	

3	./data/SRR446029_1.fastq	A1A	A1	Avr.1h.A	M1A
4	./data/SRR446030_1.fastq	A1B	A1	Avr.1h.B	M1B

## 3 Read preprocessing

### 3.1 Read quality filtering and trimming

The following example shows how one can design a custom read preprocessing function using utilities provided by the *ShortRead* package, and then apply it with `preprocessReads` in batch mode to all FASTQ samples referenced in the corresponding `SYSargs` instance (`args` object below). More detailed information on read preprocessing is provided in *systemPipeR*'s main vignette.

```
args <- systemArgs(sysma="param/trim.param", mytargets="targets_chip.txt")
filterFct <- function(fq, cutoff=20, Nexceptions=0) {
  qcount <- rowSums(as(quality(fq), "matrix") <= cutoff)
  fq[qcount <= Nexceptions] # Retains reads where Phred scores are >= cutoff with N exceptions
}
preprocessReads(args=args, Fct="filterFct(fq, cutoff=20, Nexceptions=0)", batchsize=100000)
writeTargetsout(x=args, file="targets_chip_trim.txt", overwrite=TRUE)
```

### 3.2 FASTQ quality report

The following `seeFastq` and `seeFastqPlot` functions generate and plot a series of useful quality statistics for a set of FASTQ files including per cycle quality box plots, base proportions, base-level quality trends, relative k-mer diversity, length and occurrence distribution of reads, number of reads above quality cutoffs and mean quality distribution. The results are written to a PDF file named `fastqReport.pdf`.

```
args <- systemArgs(sysma="param/bowtieSE.param", mytargets="targets_chip_trim.txt")
fqlist <- seeFastq(fastq=infile1(args), batchsize=100000, klength=8)
pdf("./results/fastqReport.pdf", height=18, width=4*length(fqlist))
seeFastqPlot(fqlist)
dev.off()
```

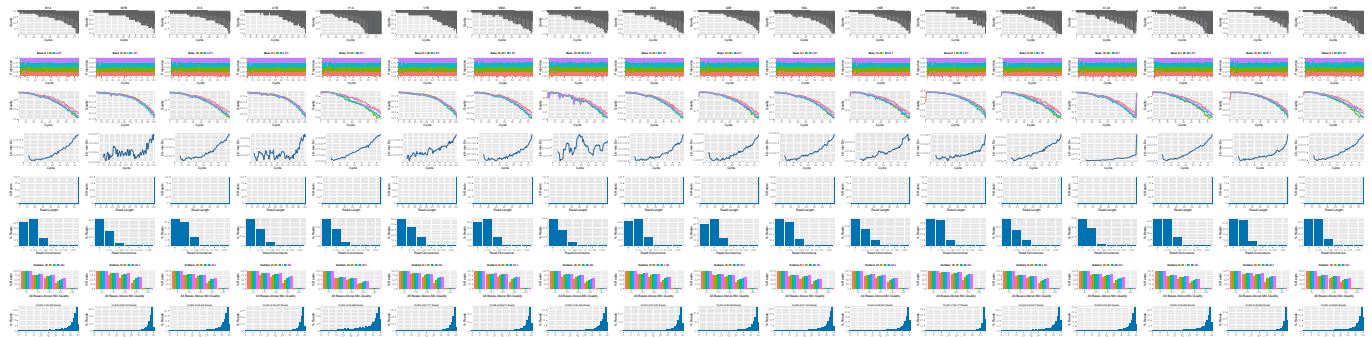


Figure 1: QC report for 18 FASTQ files.

## 4 Alignments

## 4.1 Read mapping with Bowtie2

The NGS reads of this project will be aligned with Bowtie2 against the reference genome sequence ([Langmead and Salzberg, 2012](#)). The parameter settings of the aligner are defined in the `bowtieSE.param` file. In ChIP-Seq experiments it is usually more appropriate to eliminate reads mapping to multiple locations. To achieve this, users want to remove the argument setting `'-k 50 --non-deterministic'` in the `bowtieSE.param` file.

```
args <- systemArgs(sysma="param/bowtieSE.param", mytargets="targets_chip_trim.txt")
sysargs(args)[1] # Command-line parameters for first FASTQ file
moduleload(modules(args)) # Skip if a module system is not used
system("bowtie2-build ./data/tair10.fasta ./data/tair10.fasta") # Indexes reference genome
runCommandline(args)
writeTargetsout(x=args, file="targets_bam.txt", overwrite=TRUE)
```

Check whether all BAM files have been created

```
file.exists(outpaths(args))
```

## 4.2 Read and alignment stats

The following provides an overview of the number of reads in each sample and how many of them aligned to the reference.

```
read_statsDF <- alignStats(args=args)
write.table(read_statsDF, "results/alignStats.xls", row.names=FALSE, quote=FALSE, sep="\t")
read.delim("results/alignStats.xls")
```

## 4.3 Create symbolic links for viewing BAM files in IGV

The `symLink2bam` function creates symbolic links to view the BAM alignment files in a genome browser such as IGV without moving these large files to a local system. The corresponding URLs are written to a file with a path specified under `urlfile`, here `IGVurl.txt`.

```
symLink2bam(sysargs=args, htmlDir=c("~/html/", "somedir/"),
            urlbase="http://biocluster.ucr.edu/~tgirke/",
            urlfile="./results/IGVurl.txt")
```

# 5 Peak calling with MACS2

## 5.1 Merge BAM files of replicates prior to peak calling

Merging BAM files of technical and/or biological replicates can improve the sensitivity of the peak calling by increasing the depth of read coverage. The `mergeBamByFactor` function merges BAM files based on grouping information specified by a factor, here the `Factor` column of the imported targets file. It also returns an updated `SYSargs` object containing the paths to the merged BAM files as well as to any unmerged files without replicates. This step can be skipped if merging of BAM files is not desired.

```
args <- systemArgs(sysma=NULL, mytargets="targets_bam.txt")
args_merge <- mergeBamByFactor(args, overwrite=TRUE)
writeTargetsout(x=args_merge, file="targets_mergeBamByFactor.txt", overwrite=TRUE)
```

## 5.2 Peak calling without input/reference sample

MACS2 can perform peak calling on ChIP-Seq data with and without input samples (Zhang et al., 2008). The following performs peak calling without input on all samples specified in the corresponding `args` object. Note, due to the small size of the sample data, MACS2 needs to be run here with the `'--nomodel'` setting. For real data sets, users want to remove this parameter in the corresponding `*.param` file(s).

```
args <- systemArgs(sysma="param/mac2_noinput.param", mytargets="targets_mergeBamByFactor.txt")
sysargs(args)[1] # Command-line parameters for first FASTQ file
runCommandline(args)
file.exists(outpaths(args))
writeTargetsout(x=args, file="targets_mac2.txt", overwrite=TRUE)
```

## 5.3 Peak calling with input/reference sample

To perform peak calling with input samples, they can be most conveniently specified in the `SampleReference` column of the initial targets file. The function `writeTargetsRef` uses this information to create a targets file intermediate for running MACS2 with the corresponding input samples.

```
writeTargetsRef(infile="targets_mergeBamByFactor.txt", outfile="targets_bam_ref.txt", silent=FALSE, overwr
args <- systemArgs(sysma="param/mac2.param", mytargets="targets_bam_ref.txt")
sysargs(args)[1] # Command-line parameters for first FASTQ file
runCommandline(args)
file.exists(outpaths(args))
writeTargetsout(x=args, file="targets_mac2.txt", overwrite=TRUE)
```

The peak calling results from MACS2 are written for each sample to separate files in the results directory. They are named after the corresponding `*.bam` files with extensions used by MACS2.

# 6 Annotate peaks with genomic context

## 6.1 Annotation with ChIPpeakAnno package

The following annotates the identified peaks with genomic context information using the *ChIPpeakAnno* and *ChIPseeker* packages, respectively (Zhu et al., 2010; Yu et al., 2015).

```
library(ChIPpeakAnno); library(GenomicFeatures)
args <- systemArgs(sysma="param/annotate_peaks.param", mytargets="targets_mac2.txt")
txdb <- loadDb("./data/tair10.sqlite")
ge <- genes(txdb, columns=c("tx_name", "gene_id", "tx_type"))
for(i in seq(along=args)) {
  peaksGR <- as(read.delim(infile1(args)[i], comment="#"), "GRanges")
  annotatedPeak <- annotatePeakInBatch(peaksGR, AnnotationData=genes(txdb))
  df <- data.frame(as.data.frame(annotatedPeak), as.data.frame(values(ge[values(annotatedPeak)$feature,])
  write.table(df, outpaths(args[i]), quote=FALSE, row.names=FALSE, sep="\t")
}
writeTargetsout(x=args, file="targets_peakanno.txt", overwrite=TRUE)
```

The peak annotation results are written for each peak set to separate files in the results directory. They are named after the corresponding peak files with extensions specified in the `annotate_peaks.param` file, here `*.peaks.annotated.xls`.

## 6.2 Annotation with ChIPseeker package

Same as in previous step but using the *ChIPseeker* package for annotating the peaks.

```
library(ChIPseeker)
txdb <- loadDb("./data/tair10.sqlite")
for(i in seq(along=args)) {
  peakAnno <- annotatePeak(infile1(args)[i], TxDb=txdb, verbose=FALSE)
  df <- as.data.frame(peakAnno)
  write.table(df, outpaths(args[i]), quote=FALSE, row.names=FALSE, sep="\t")
}
writeTargetsout(x=args, file="targets_peakanno.txt", overwrite=TRUE)
```

Summary plots provided by the *ChIPseeker* package. Here applied only to one sample for demonstration purposes.

```
peak <- readPeakFile(infile1(args)[1])
covplot(peak, weightCol="X.log10.pvalue.")
peakHeatmap(outpaths(args)[1], TxDb=txdb, upstream=1000, downstream=1000, color="red")
plotAvgProf2(outpaths(args)[1], TxDb=txdb, upstream=1000, downstream=1000, xlab="Genomic Region (5'→3')",
```

## 7 Count reads overlapping peak regions

The `countRangeset` function is a convenience wrapper to perform read counting iteratively over several range sets, here peak range sets. Internally, the read counting is performed with the `summarizeOverlaps` function from the *GenomicAlignments* package. The resulting count tables are directly saved to files, one for each peak set.

```
library(GenomicRanges)
args <- systemArgs(sysma="param/count_rangesets.param", mytargets="targets_macs.txt")
args_bam <- systemArgs(sysma=NULL, mytargets="targets_bam.txt")
bfl <- BamFileList(outpaths(args_bam), yieldSize=50000, index=character())
countDFnames <- countRangeset(bfl, args, mode="Union", ignore.strand=TRUE)
writeTargetsout(x=args, file="targets_countDF.txt", overwrite=TRUE)
```

## 8 Differential binding analysis of peaks

The function `runDiff` performs differential binding analysis in batch mode for several count tables using *edgeR* or *DESeq2* (Robinson et al., 2010; Love et al., 2014). Internally, it calls the functions `run_edgeR` and `run_DESeq2`. It also returns the filtering results and plots from the downstream `filterDEGs` function using the fold change and FDR cutoffs provided under the `dbrfilter` argument.

```
args_diff <- systemArgs(sysma="param/rundiff.param", mytargets="targets_countDF.txt")
cmp <- readComp(file=args_bam, format="matrix")
dbrlist <- runDiff(args=args_diff, diffFct=run_edgeR, targets=targetsin(args_bam),
  cmp=cmp[[1]], independent=TRUE, dbrfilter=c(Fold=2, FDR=1))
writeTargetsout(x=args_diff, file="targets_rundiff.txt", overwrite=TRUE)
```

## 9 GO term enrichment analysis

The following performs GO term enrichment analysis for each annotated peak set.

```
args <- systemArgs(sysma="param/mac2.param", mytargets="targets_bam_ref.txt")
args_anno <- systemArgs(sysma="param/annotate_peaks.param", mytargets="targets_mac2.txt")
annofiles <- outpaths(args_anno)
gene_ids <- sapply(names(annofiles), function(x) unique(as.character(read.delim(annofiles[x])[, "gene_id"])))
load("data/G0/catdb.RData")
BatchResult <- GOCluster_Report(catdb=catdb, setlist=gene_ids, method="all", id_type="gene", CLSZ=2, cutoff=0.05)
```

## 10 Motif analysis

### 10.1 Parse DNA sequences of peak regions from genome

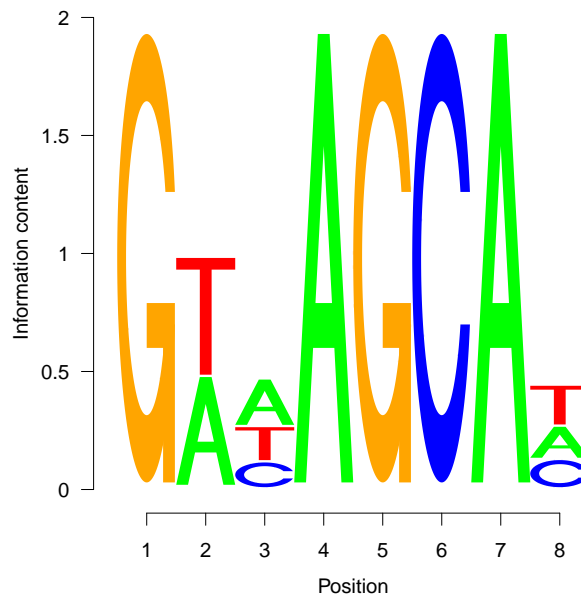
Enrichment analysis of known DNA binding motifs or *de novo* discovery of novel motifs requires the DNA sequences of the identified peak regions. To parse the corresponding sequences from the reference genome, the `getSeq` function from the *Biostrings* package can be used. The following example parses the sequences for each peak set and saves the results to separate FASTA files, one for each peak set. In addition, the sequences in the FASTA files are ranked (sorted) by increasing p-values as expected by some motif discovery tools, such as *BCRANK*.

```
library(Biostrings); library(seqLogo); library(BCRANK)
args <- systemArgs(sysma="param/annotate_peaks.param", mytargets="targets_mac2.txt")
rangefiles <- infile1(args)
for(i in seq(along=rangefiles)) {
  df <- read.delim(rangefiles[i], comment="#")
  peaks <- as(df, "GRanges")
  names(peaks) <- paste0(as.character(seqnames(peaks)), "_", start(peaks), "-", end(peaks))
  peaks <- peaks[order(values(peaks)$X.log10.pvalue, decreasing=TRUE)]
  pseq <- getSeq(FaFile("./data/tair10.fasta"), peaks)
  names(pseq) <- names(peaks)
  writeXStringSet(pseq, paste0(rangefiles[i], ".fasta"))
}
```

### 10.2 Motif discovery with BCRANK

The Bioconductor package *BCRANK* is one of the many tools available for *de novo* discovery of DNA binding motifs in peak regions of ChIP-Seq experiments. The given example applies this method on the first peak sample set and plots the sequence logo of the highest ranking motif.

```
set.seed(0)
BCRANKout <- bcrank(paste0(rangefiles[1], ".fasta"), restarts=25, use.P1=TRUE, use.P2=TRUE)
topTable(BCRANKout)
topMotif <- topTable(BCRANKout, 1)
weightMatrix <- pwm(topMotif, normalize = FALSE)
weightMatrixNormalized <- pwm(topMotif, normalize = TRUE)
pdf("results/seqlogo.pdf")
seqLogo(weightMatrixNormalized)
dev.off()
```

Figure 2: One of the motifs identified by *BCRANK*

## 11 Version Information

```
toLatex(sessionInfo())
```

- R version 3.2.5 (2016-04-14), x86\_64-pc-linux-gnu
- Locale: LC\_CTYPE=en\_US.UTF-8, LC\_NUMERIC=C, LC\_TIME=en\_US.UTF-8, LC\_COLLATE=en\_US.UTF-8, LC\_MONETARY=en\_US.UTF-8, LC\_MESSAGES=en\_US.UTF-8, LC\_PAPER=en\_US.UTF-8, LC\_NAME=C, LC\_ADDRESS=C, LC\_TELEPHONE=C, LC\_MEASUREMENT=en\_US.UTF-8, LC\_IDENTIFICATION=C
- Base packages: base, datasets, graphics, grDevices, methods, parallel, stats, stats4, utils
- Other packages: Biobase 2.30.0, BiocGenerics 0.16.1, BiocParallel 1.4.0, Biostrings 2.38.2, DBI 0.3.1, GenomInfoDb 1.6.1, GenomicAlignments 1.6.1, GenomicRanges 1.22.1, IRanges 2.4.4, knitr 1.12, Rsamtools 1.22.0, RSQLite 1.0.0, S4Vectors 0.8.3, ShortRead 1.28.0, SummarizedExperiment 1.0.1, systemPipeR 1.4.8, XVector 0.10.0
- Loaded via a namespace (and not attached): annotate 1.48.0, AnnotationDbi 1.32.0, AnnotationForge 1.12.0, base64enc 0.1-3, BatchJobs 1.6, BBmisc 1.9, BiocStyle 1.8.0, biomaRt 2.26.1, bitops 1.0-6, brew 1.0-6, Category 2.36.0, checkmate 1.6.3, colorspace 1.2-6, digest 0.6.9, edgeR 3.12.0, evaluate 0.8, fail 1.3, formatR 1.2.1, futile.logger 1.4.1, futile.options 1.0.0, genefilter 1.52.0, GenomicFeatures 1.22.6, ggplot2 2.0.0, GO.db 3.2.2, GOSTats 2.36.0, graph 1.48.0, grid 3.2.5, GSEABase 1.32.0, gtable 0.1.2, highr 0.5.1, hwriter 1.3.2, lambda.r 1.1.7, lattice 0.20-33, latticeExtra 0.6-26, limma 3.26.3, magrittr 1.5, Matrix 1.2-3, munsell 0.4.2, pheatmap 1.0.7, plyr 1.8.3, RBGL 1.46.0, RColorBrewer 1.1-2, Rcpp 0.12.3, RCurl 1.95-4.7, rjson 0.2.15, rtracklayer 1.30.1, scales 0.3.0, sendmailR 1.2-1, splines 3.2.5, stringi 1.0-1, stringr 1.0.0, survival 2.38-3, tools 3.2.5, XML 3.98-1.3, xtable 1.8-0, zlibbioc 1.16.0



## 12 Funding

---

This project was supported by funds from the National Institutes of Health (NIH) and the National Science Foundation (NSF).

## 13 References

---

- Thomas Girke. systemPipeR: NGS workflow and report generation environment, 28 June 2014. URL <https://github.com/tgirke/systemPipeR>.
- Ben Langmead and Steven L Salzberg. Fast gapped-read alignment with bowtie 2. *Nat. Methods*, 9(4):357–359, April 2012. ISSN 1548-7091. doi: 10.1038/nmeth.1923. URL <http://dx.doi.org/10.1038/nmeth.1923>.
- Michael Love, Wolfgang Huber, and Simon Anders. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.*, 15(12):550, 2014. ISSN 1465-6906. doi: 10.1186/s13059-014-0550-8. URL <http://genomebiology.com/2014/15/12/550>.
- M D Robinson, D J McCarthy, and G K Smyth. edgeR: a bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*, 26(1):139–140, January 2010. ISSN 1367-4803. doi: 10.1093/bioinformatics/btp616. URL <http://dx.doi.org/10.1093/bioinformatics/btp616>.
- Guangchuang Yu, Li-Gen Wang, and Qing-Yu He. ChIPseeker: an R/Bioconductor package for ChIP peak annotation, comparison and visualization. *Bioinformatics*, 31(14):2382–2383, 15 July 2015. ISSN 1367-4803, 1367-4811. doi: 10.1093/bioinformatics/btv145. URL <http://dx.doi.org/10.1093/bioinformatics/btv145>.
- Y Zhang, T Liu, C A Meyer, J Eeckhoute, D S Johnson, B E Bernstein, C Nussbaum, R M Myers, M Brown, W Li, and X S Liu. Model-based analysis of ChIP-Seq (MACS). *Genome Biol.*, 9(9), 2008. ISSN 1465-6906. doi: 10.1186/gb-2008-9-9-r137. URL <http://dx.doi.org/10.1186/gb-2008-9-9-r137>.
- Lihua J Zhu, Claude Gazin, Nathan D Lawson, Hervé Pagès, Simon M Lin, David S Lapointe, and Michael R Green. ChIPpeakAnno: a bioconductor package to annotate ChIP-seq and ChIP-chip data. *BMC Bioinformatics*, 11:237, 11 May 2010. ISSN 1471-2105. doi: 10.1186/1471-2105-11-237. URL <http://dx.doi.org/10.1186/1471-2105-11-237>.