

# R in NGS 实验 3

生信 2001 张子栋 2020317210101

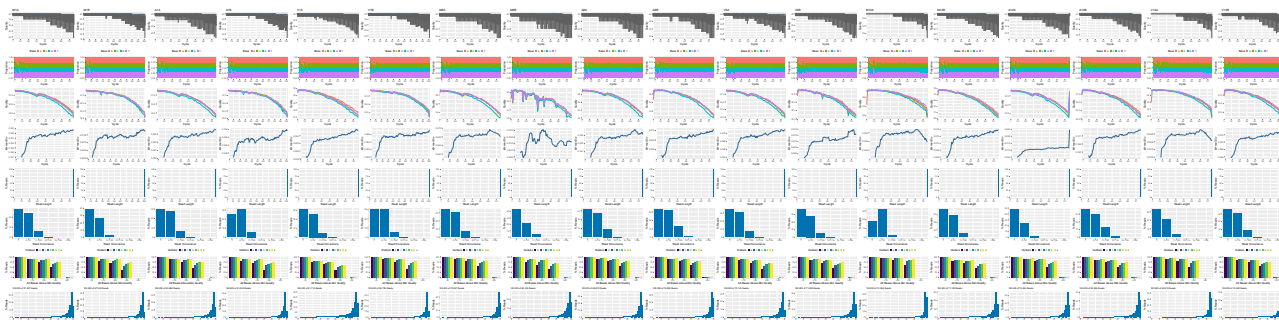
## 1 实验步骤及结果

```
1 > library(systemPipeR)
2 > library(systemPipeRdata)
3 > setwd(choose.dir())
4 > genWorkenvir(workflow = "chipseq")
5 > targetsPath <- system.file("extdata", "targets_chip.txt", package =
  "systemPipeR")
6 > targets <- read.delim(targetsPath, comment.char = "#")
7 > targets[1:4, -c(5, 6)]
8
      FileName SampleName Factor SampleLong SampleReference
9 1 ./data/SRR446027_1.fastq.gz      M1A      M1 Mock.1h.A
10 2 ./data/SRR446028_1.fastq.gz      M1B      M1 Mock.1h.B
11 3 ./data/SRR446029_1.fastq.gz      A1A      A1 Avr.1h.A      M1A
12 4 ./data/SRR446030_1.fastq.gz      A1B      A1 Avr.1h.B      M1B
13 > dir_path <- system.file("extdata/cwl/preprocessReads/trim-se", package =
  "systemPipeR")
14 > trim <- loadWF(targets = targetsPath, wf_file = "trim-se.cwl", input_file =
  "trim-se.yml", dir_path = dir_path)
15 > trim <- renderWF(trim, inputvars = c(FileName = "_FASTQ_PATH1_", SampleName =
  "_SampleName_"))
16 > output(trim)[1:2]
17 $M1A
18 $M1A$`trim-se`
19 [1] "./results/M1A.fastq_trim.gz"
20
21
22 $M1B
23 $M1B$`trim-se`
24 [1] "./results/M1B.fastq_trim.gz"
25 > filterFct <- function(fq, cutoff = 20, Nexceptions = 0) {
26 +   qcount <- rowSums(as(quality(fq), "matrix") <= cutoff, na.rm = TRUE)
27 +   fq[qcount <= Nexceptions]
28 +   # Retains reads where Phred scores are >= cutoff with N
29 +   # exceptions
30 + }
31 > preprocessReads(args = trim, Fct = "filterFct(fq, cutoff=20, Nexceptions=0)",
32 +   batchSize = 1e+05)
33 44022 processed reads written to file: ./results/M1A.fastq_trim.gz
34 44927 processed reads written to file: ./results/M1B.fastq_trim.gz
35 47793 processed reads written to file: ./results/A1A.fastq_trim.gz
36 41201 processed reads written to file: ./results/A1B.fastq_trim.gz
37 38549 processed reads written to file: ./results/V1A.fastq_trim.gz
```

```

38 49362 processed reads written to file: ./results/V1B.fastq_trim.gz
39 58018 processed reads written to file: ./results/M6A.fastq_trim.gz
40 41708 processed reads written to file: ./results/M6B.fastq_trim.gz
41 54241 processed reads written to file: ./results/A6A.fastq_trim.gz
42 62722 processed reads written to file: ./results/A6B.fastq_trim.gz
43 52165 processed reads written to file: ./results/V6A.fastq_trim.gz
44 50684 processed reads written to file: ./results/V6B.fastq_trim.gz
45 41583 processed reads written to file: ./results/M12A.fastq_trim.gz
46 49316 processed reads written to file: ./results/M12B.fastq_trim.gz
47 56692 processed reads written to file: ./results/A12A.fastq_trim.gz
48 58547 processed reads written to file: ./results/A12B.fastq_trim.gz
49 42475 processed reads written to file: ./results/V12A.fastq_trim.gz
50 56732 processed reads written to file: ./results/V12B.fastq_trim.gz
51 > writeTargetsout(x = trim, file = "targets_chip_trim.txt", step = 1,
52 +               new_col = "FileName", new_col_output_index = 1, overwrite =
    TRUE)
53     Written content of 'targetsout(x)' to file: targets_chip_trim.txt
54 > library(BiocParallel)
55 > library(batchtools)
56 > f <- function(x) {
57 +   targets <- system.file("extdata", "targets_chip.txt", package =
    "systemPipeR")
58 +   dir_path <- system.file("extdata/cwl/preprocessReads/trim-se",
    package = "systemPipeR")
59 +   trim <- loadWorkflow(targets = targets, wf_file = "trim-se.cwl",
60 +   input_file = "trim-se.yml", dir_path = dir_path)
61 +   trim <- renderWF(trim, inputvars = c(FileName = "_FASTQ_PATH1_",
62 +   SampleName = "_SampleName_"))
63 +   seeFastq(fastq = infile1(trim)[x], batchsize = 1e+05, klength = 8)
64 + }
65 > resources <- list(walltime = 120, ntasks = 1, ncpus = 4, memory = 1024)
66 > fqlist <- lapply(seq(along = trim), f)
67 > pdf("./results/fastqReport.pdf", height = 18, width = 4 * length(fqlist))

```



矢量图，可放大。

```

1 > seeFastqPlot(unlist(fqlist, recursive = FALSE))
2 > dev.off()
3 pdf

```

```

4   2
5   > library(ChIPpeakAnno)
6   > library(GenomicFeatures)
7   > dir_path <- system.file("extdata/cwl/annotate_peaks", package = "systemPipeR")
8   > dir_path <- system.file("extdata/cwl/annotate_peaks", package = "systemPipeR")
9   > args <- loadWF(targets = "targets_macs.txt", wf_file = "annotate_peaks.cwl",
10  +               input_file = "annotate_peaks.yml", dir_path = dir_path)
11  > args <- renderWF(args, inputvars = c(FileName = "_FASTQ_PATH1_",
12  +                                     SampleName = "_SampleName_"))
13  >
14  > txdb <- makeTxDbFromGFF(file = "data/tair10.gff", format = "gff",
15  +                       dataSource = "TAIR", organism = "Arabidopsis thaliana")
16  > for (i in seq(along = args)) {
17  +   peaksGR <- as(read.delim(infile1(args)[i], comment = "#"),
18  +                 "GRanges")
19  +   annotatedPeak <- annotatePeakInBatch(peaksGR, AnnotationData = genes(txdb))
20  +   df <- data.frame(as.data.frame(annotatedPeak),
21  +                   as.data.frame(values(ge[values(annotatedPeak)$feature,
22  +                                     ])))
23  +   df$tx_type <- unlist(df$tx_type)
24  +   tx_name <- c()
25  +   for (j in df$tx_name){
26  +     tx_name <- rbind(tx_name, j[1])
27  +   }
28  +   df$tx_name <- tx_name
29  +   outpaths <- subsetWF(args, slot = "output", subset = 1, index = 1)
30  +   write.table(df, outpaths[i], quote = FALSE, row.names = FALSE,
31  +               sep = "\t")
32  + }
33  > writeTargetsout(x = args, file = "targets_peakanno.txt", step = 1,
34  +                 new_col = "FileName", new_col_output_index = 1, overwrite =
35  +                 TRUE)
36  > library(ChIPseeker)
37  > for (i in seq(along = args)) {
38  +   peakAnno <- annotatePeak(infile1(args)[i], TxDb = txdb, verbose = FALSE)
39  +   df <- as.data.frame(peakAnno)
40  +   outpaths <- subsetWF(args, slot = "output", subset = 1, index = 1)
41  +   write.table(df, outpaths[i], quote = FALSE, row.names = FALSE,
42  +               sep = "\t")
43  + }
44  > writeTargetsout(x = args, file = "targets_peakanno.txt", step = 1,
45  +                 new_col = "FileName", new_col_output_index = 1, overwrite =
46  +                 TRUE)
47  > library(ChIPseeker)
48  > for (i in seq(along = args)) {
49  +   peakAnno <- annotatePeak(infile1(args)[i], TxDb = txdb, verbose = FALSE)
50  +   df <- as.data.frame(peakAnno)

```



```

15 > outpaths <- subsetWF(args_bam, slot = "output", subset = 1, index = 1)
16 > bfl <- BamFileList(outpaths, yieldSize = 50000, index = character())
17 > # countDFnames <- countRangeset(bfl, args, mode = "Union", ignore.strand =
TRUE) # skipped
18 > # writeTargetsout(x = args, file = "targets_countDF.txt", step = 1, new_col =
"FileName", new_col_output_index = 1, overwrite = TRUE) # skipped
19 > library(GenomicRanges)
20 > dir_path <- system.file("extdata/cwl/count_rangesets", package = "systemPipeR")
21 > args <- loadWF(targets = "targets_mac3.txt", wf_file = "count_rangesets.cwl",
22 +             input_file = "count_rangesets.yml", dir_path = dir_path)
23 > args <- renderWF(args, inputvars = c(FileName = "_FASTQ_PATH1_",
24 +             SampleName = "_SampleName_"))
25 >
26 > ## Bam Files
27 > targets <- system.file("extdata", "targets_chip.txt", package = "systemPipeR")
28 > dir_path <- system.file("extdata/cwl/bowtie2", package = "systemPipeR")
29 > args_bam <- loadWF(targets = targets, wf_file = "bowtie2-mapping-se.cwl",
30 +             input_file = "bowtie2-mapping-se.yml", dir_path = dir_path)
31 > args_bam <- renderWF(args_bam, inputvars = c(FileName = "_FASTQ_PATH1_",
32 +             SampleName = "_SampleName_"))
33 > args_bam <- output_update(args_bam, dir = FALSE, replace = TRUE,
34 +             extension = c(".sam", ".bam"))
35 > outpaths <- subsetWF(args_bam, slot = "output", subset = 1, index = 1)
36 > bfl <- BamFileList(outpaths, yieldSize = 50000, index = character())
37 > countDFnames <- countRangeset(bfl, args, mode = "Union", ignore.strand = TRUE)
38 > dir_path <- system.file("extdata/cwl/rundiff", package = "systemPipeR")
39 > args_diff <- loadWF(targets = "targets_countDF.txt", wf_file = "rundiff.cwl",
40 +             input_file = "rundiff.yml", dir_path = dir_path)
41 > args_diff <- renderWF(args_diff, inputvars = c(FileName = "_FASTQ_PATH1_",
42 +             SampleName = "_SampleName_"))
43 >
44 > cmp <- readComp(file = args_bam, format = "matrix")
45 > dbrlist <- runDiff(args = args_diff, diffFct = run_edgeR, targets =
targets.as.df(targets(args_bam)),
46 +             cmp = cmp[[1]], independent = TRUE, dbrfilter = c(Fold = 2,
47 +             FDR = 1))
48 Disp = 0.25073 , BCV = 0.5007
49 Disp = 0.15223 , BCV = 0.3902
50 Disp = 0.19792 , BCV = 0.4449
51 Disp = 0.11692 , BCV = 0.3419
52 Disp = 0.09286 , BCV = 0.3047
53 Disp = 0.14312 , BCV = 0.3783
54 Disp = 0.17049 , BCV = 0.4129
55 Disp = 0.09671 , BCV = 0.311
56 Disp = 0.14805 , BCV = 0.3848
57 Wrote count result 1 to M1A_peaks.edgeR.xls
58 Saved plot 1 to M1A_peaks.edgeR.xls.pdf

```

```
59 Disp = 0.23843 , BCV = 0.4883
60 Disp = 0.11547 , BCV = 0.3398
61 Disp = 0.21655 , BCV = 0.4653
62 Disp = 0.04269 , BCV = 0.2066
63 Disp = 0.15216 , BCV = 0.3901
64 Disp = 0.12235 , BCV = 0.3498
65 Disp = 0.20886 , BCV = 0.457
66 Disp = 0.12322 , BCV = 0.351
67 Disp = 0.20853 , BCV = 0.4567
68 Wrote count result 2 to A1A_peaks.edgeR.xls
69 Saved plot 2 to A1A_peaks.edgeR.xls.pdf
70 Disp = 0.10458 , BCV = 0.3234
71 Disp = 0.13358 , BCV = 0.3655
72 Disp = 0.09085 , BCV = 0.3014
73 Disp = 0.06552 , BCV = 0.256
74 Disp = 0.13023 , BCV = 0.3609
75 Disp = 0.1174 , BCV = 0.3426
76 Disp = 0.09355 , BCV = 0.3059
77 Disp = 0.08648 , BCV = 0.2941
78 Disp = 0.0827 , BCV = 0.2876
79 Wrote count result 3 to V1A_peaks.edgeR.xls
80 Saved plot 3 to V1A_peaks.edgeR.xls.pdf
81 Disp = 0.19174 , BCV = 0.4379
82 Disp = 0.12306 , BCV = 0.3508
83 Disp = 0.15138 , BCV = 0.3891
84 Disp = 0.13047 , BCV = 0.3612
85 Disp = 0.09703 , BCV = 0.3115
86 Disp = 0.16206 , BCV = 0.4026
87 Disp = 0.17529 , BCV = 0.4187
88 Disp = 0.1089 , BCV = 0.33
89 Disp = 0.15879 , BCV = 0.3985
90 Wrote count result 4 to M6A_peaks.edgeR.xls
91 Saved plot 4 to M6A_peaks.edgeR.xls.pdf
92 Disp = 0.10874 , BCV = 0.3298
93 Disp = 0.17742 , BCV = 0.4212
94 Disp = 0.12824 , BCV = 0.3581
95 Disp = 0.24666 , BCV = 0.4967
96 Disp = 0.10608 , BCV = 0.3257
97 Disp = 0.20981 , BCV = 0.4581
98 Disp = 0.23416 , BCV = 0.4839
99 Disp = 0.10848 , BCV = 0.3294
100 Disp = 0.21298 , BCV = 0.4615
101 Wrote count result 5 to A6A_peaks.edgeR.xls
102 Saved plot 5 to A6A_peaks.edgeR.xls.pdf
103 Disp = 0.12458 , BCV = 0.353
104 Disp = 0.12455 , BCV = 0.3529
105 Disp = 0.1462 , BCV = 0.3824
```

```
106 Disp = 0.0498 , BCV = 0.2232
107 Disp = 0.05905 , BCV = 0.243
108 Disp = 0.06298 , BCV = 0.251
109 Disp = 0.10963 , BCV = 0.3311
110 Disp = 0.09929 , BCV = 0.3151
111 Disp = 0.10645 , BCV = 0.3263
112 Wrote count result 6 to V6A_peaks.edgeR.xls
113 Saved plot 6 to V6A_peaks.edgeR.xls.pdf
114 Disp = 0.19246 , BCV = 0.4387
115 Disp = 0.10992 , BCV = 0.3315
116 Disp = 0.15488 , BCV = 0.3935
117 Disp = 0.12168 , BCV = 0.3488
118 Disp = 0.09003 , BCV = 0.3
119 Disp = 0.14578 , BCV = 0.3818
120 Disp = 0.16995 , BCV = 0.4122
121 Disp = 0.09506 , BCV = 0.3083
122 Disp = 0.15563 , BCV = 0.3945
123 Wrote count result 7 to M12A_peaks.edgeR.xls
124 Saved plot 7 to M12A_peaks.edgeR.xls.pdf
125 Disp = 0.29779 , BCV = 0.5457
126 Disp = 0.16126 , BCV = 0.4016
127 Disp = 0.26423 , BCV = 0.514
128 Disp = 0.28419 , BCV = 0.5331
129 Disp = 0.11531 , BCV = 0.3396
130 Disp = 0.28049 , BCV = 0.5296
131 Disp = 0.27675 , BCV = 0.5261
132 Disp = 0.13993 , BCV = 0.3741
133 Disp = 0.25 , BCV = 0.5
134 Wrote count result 8 to A12A_peaks.edgeR.xls
135 Saved plot 8 to A12A_peaks.edgeR.xls.pdf
136 Disp = 0.22277 , BCV = 0.472
137 Disp = 0.19447 , BCV = 0.441
138 Disp = 0.17158 , BCV = 0.4142
139 Disp = 0.23746 , BCV = 0.4873
140 Disp = 0.14669 , BCV = 0.383
141 Disp = 0.24076 , BCV = 0.4907
142 Disp = 0.20954 , BCV = 0.4578
143 Disp = 0.11004 , BCV = 0.3317
144 Disp = 0.21206 , BCV = 0.4605
145 Wrote count result 9 to V12A_peaks.edgeR.xls
146 Saved plot 9 to V12A_peaks.edgeR.xls.pdf
147 > writeTargetsout(x = args_diff, file = "targets_rundiff.txt",
148 +                 step = 1, new_col = "FileName", new_col_output_index = 1,
149 +                 overwrite = TRUE)
150     Written content of 'targetsout(x)' to file: targets_rundiff.txt
151 > writeTargetsout(x = args_diff, file = "targets_rundiff.txt",
152 +                 step = 1, new_col = "FileName", new_col_output_index = 1,
```

```

153 +             overwrite = TRUE)
154     Written content of 'targetsout(x)' to file: targets_rundiff.txt
155 > dir_path <- system.file("extdata/cwl/annotate_peaks", package = "systemPipeR")
156 > args <- loadWF(targets = "targets_bam_ref.txt", wf_file = "annotate_peaks.cwl",
157 +             input_file = "annotate_peaks.yml", dir_path = dir_path)
158 > args <- renderWF(args, inputvars = c(FileName1 = "_FASTQ_PATH1_",
159 +             SampleName = "_SampleName_"))
160 >
161 > args_anno <- loadWF(targets = "targets_macs.txt", wf_file =
162 +             "annotate_peaks.cwl",
163 +             input_file = "annotate_peaks.yml", dir_path = dir_path)
164 > args_anno <- renderWF(args_anno, inputvars = c(FileName = "_FASTQ_PATH1_",
165 +             SampleName = "_SampleName_"))
166 > annofiles <- subsetWF(args_anno, slot = "output", subset = 1,
167 +             index = 1)
168 > gene_ids <- sapply(names(annofiles), function(x)
169 +             unique(as.character(read.delim(annofiles[x])[,
170 +             "geneId"]))), simplify = FALSE)
171 > load("data/GO/catdb.RData")
172 > BatchResult <- GOCluster_Report(catdb = catdb, setlist = gene_ids,
173 +             method = "all", id_type = "gene", CLSZ = 2,
174 +             cutoff = 0.9,
175 +             gocats = c("MF", "BP", "CC"), recordSpecGO =
176 +             NULL)
177 > library(Biostrings)
178 > library(seqLogo)
179 > library(BCRANK)
180 > dir_path <- system.file("extdata/cwl/annotate_peaks", package = "systemPipeR")
181 > args <- loadWF(targets = "targets_macs.txt", wf_file = "annotate_peaks.cwl",
182 +             input_file = "annotate_peaks.yml", dir_path = dir_path)
183 > args <- renderWF(args, inputvars = c(FileName = "_FASTQ_PATH1_",
184 +             SampleName = "_SampleName_"))
185 >
186 > rangefiles <- infile1(args)
187 > for (i in seq(along = rangefiles)) {
188 +     df <- read.delim(rangefiles[i], comment = "#")
189 +     peaks <- as(df, "GRanges")
190 +     names(peaks) <- paste0(as.character(seqnames(peaks)), "_",
191 +             start(peaks), "-", end(peaks))
192 +     peaks <- peaks[order(values(peaks)$X.log10.pvalue., decreasing = TRUE)]
193 +     pseq <- getSeq(FaFile("./data/tair10.fasta"), peaks)
194 +     names(pseq) <- names(peaks)
195 +     writeXStringSet(pseq, paste0(rangefiles[i], ".fasta"))
196 + }
197 > set.seed(0)
198 > BCRANKout <- bcrank(paste0(rangefiles[1], ".fasta"), restarts = 25,

```



```

195 + use.P1 = TRUE, use.P2 = TRUE)
196 > toptable(BCRANKout)
197     Consensus    Score
198 1    AGTCAHTT 87.08268
199 2    ATGTNAGA 78.86206
200 3    CACAHDAA 78.43056
201 4    MGGTATC 77.63990
202 5    AMRCABAAR 69.68058
203 6    AWAARBCAA 69.00625
204 7    DCCDDGAAAS 62.46803
205 8    RYNTGNCTCT 61.84035
206 9    DHGABWGGAA 61.19199
207 10   TGNSTTCHT 60.77090
208 11   GHTGABNTTM 60.46803
209 12   GCHGHTVTVT 58.78979
210 13   CABKBTGBHA 58.76985
211 14   AAHTBHCTVC 57.98167
212 15   DDBCGBCCDT 57.18527
213 16   TTVGMNAGWTC 56.98050
214 17   GATTKGBNGAA 55.79763
215 18   DDHCNGHCTTG 52.71424
216 19   AKAGAAHAGC 52.56690
217 20   TBRTAKCTNC 52.18433
218 21   ACGAWNTBRK 51.60559
219 22   VWC GTVNDVTT 50.63590
220 23   CVDGDTCAVVC 50.49847
221 24   GAHNAMASAC 49.08524
222 25   CCBASGYNDG 46.22657
223 > topMotif <- toptable(BCRANKout, 1)
224 > weightMatrix <- pwm(topMotif, normalize = FALSE)
225 > weightMatrixNormalized <- pwm(topMotif, normalize = TRUE)
226 > pdf("results/seqlogo.pdf")
227 > seqLogo(weightMatrixNormalized)
228 > dev.off()
229 RStudioGD
230     2

```

