

The *exomePeak2* user's guide

| Zhen Wei <ZhenWei@xjtlu.edu.cn> | Jia Meng <JiaMeng@xjtlu.edu.cn> | Department of Biological Sciences, Xi'an Jiaotong-Liverpool University, Suzhou, Jiangsu, 215123, China | Institute of Integrative Biology, University of Liverpool, L7 8TX, Liverpool, United Kingdom

2019-11-25

Contents

1	Peak Calling	2
2	Differential Modification Analysis	3
3	Quantification and Statistical Analysis with Single Based Modification Annotation	4
4	Peak Calling and Visualization in Multiple Steps	4
5	Contact	6
6	Session Info	7

1 Peak Calling

Users need to specify the bam file directories of IP and input samples separately using the arguments of `bam_ip` and `bam_input`; the biological replicates are represented by a character vector of the **BAM** file directories.

Transcript annotation is provided using GFF files in this example. The transcript annotation can also come from the `TxDb` object. *exomePeak2* will automatically download the `TxDb` if you fill the `genome` argument with the UCSC genome name.

The genome sequence is required to conduct GC content bias correction. If the `genome` argument is missing (= `NULL`), *exomePeak2* will perform peak calling without correcting the GC content bias.

```
library(exomePeak2)

GENE_ANNO_GTF = system.file("extdata", "example.gtf", package="exomePeak2")

f1 = system.file("extdata", "IP1.bam", package="exomePeak2")
f2 = system.file("extdata", "IP2.bam", package="exomePeak2")
f3 = system.file("extdata", "IP3.bam", package="exomePeak2")
f4 = system.file("extdata", "IP4.bam", package="exomePeak2")
IP_BAM = c(f1,f2,f3,f4)

f1 = system.file("extdata", "Input1.bam", package="exomePeak2")
f2 = system.file("extdata", "Input2.bam", package="exomePeak2")
f3 = system.file("extdata", "Input3.bam", package="exomePeak2")
INPUT_BAM = c(f1,f2,f3)

exomePeak2(bam_ip = IP_BAM,
           bam_input = INPUT_BAM,
           gff_dir = GENE_ANNO_GTF,
           genome = "hg19",
           paired_end = FALSE)
## class: SummarizedExomePeak
## dim: 31 7
## metadata(0):
## assays(2): counts GCsizeFactors
## rownames(31): mod_11 mod_13 ... control_13 control_14
## rowData names(2): GC_content feature_length
## colnames(7): IP1.bam IP2.bam ... Input2.bam Input3.bam
## colData names(3): design_IP design_Treatment sizeFactor
```

exomePeak2 will export the modification peaks in formats of **BED** file and **CSV** table, the data will be saved automatically under a folder named by `exomePeak2_output`.

Explanation over the columns of the exported table:

- **chr** : the chromosomal name of the peak.
- **chromStart** : the start of the peak on the chromosome.
- **chromEnd** : the end of the peak on the chromosome.
- **name** : the unique ID of the modification peak.
- **score** : the $-\log_2$ p value of the peak.
- **strand** : the strand of the peak on genome.

The *exomePeak2* user's guide

- ***thickStart*** : the start position of the peak.
- ***thickEnd*** : the end position of the peak.
- ***itemRgb*** : the column for the RGB encoded color in BED file visualization.
- ***blockCount*** : the block (exon) number within the peak.
- ***blockSizes*** : the widths of blocks.
- ***blockStarts*** : the start positions of blocks.
- ***geneID*** : the gene ID of the peak.
- ***ReadsCount.input*** : the reads count of the input sample.
- ***ReadsCount.IP*** : the reads count of the IP sample.
- ***log2FoldChange*** : the log2 IP over input fold enrichment.
- ***pvalue*** : the p value of the enrichment.
- ***padj*** : the adjusted p value using BH approach.

2 Differential Modification Analysis

The code below could conduct differential modification analysis (Comparison of Two Conditions) on exon regions defined by the transcript annotation.

In differential modification mode, *exomePeak2* will first perform Peak calling on exon regions using both the control and treated samples. Then, it will conduct the differential modification analysis on peaks reported from peak calling using an interactive GLM.

```
f1 = system.file("extdata", "treated_IP1.bam", package="exomePeak2")
TREATED_IP_BAM = c(f1)
f1 = system.file("extdata", "treated_Input1.bam", package="exomePeak2")
TREATED_INPUT_BAM = c(f1)

exomePeak2(bam_ip = IP_BAM,
            bam_input = INPUT_BAM,
            bam_treated_input = TREATED_INPUT_BAM,
            bam_treated_ip = TREATED_IP_BAM,
            gff_dir = GENE_ANNO_GTF,
            genome = "hg19",
            paired_end = FALSE)

## class: SummarizedExomePeak
## dim: 23 9
## metadata(0):
## assays(2): counts GCsizeFactors
## rownames(23): mod_10 mod_11 ... control_5 control_6
## rowData names(2): GC_content feature_length
## colnames(9): IP1.bam IP2.bam ... treated_IP1.bam treated_Input1.bam
## colData names(3): design_IP design_Treatment sizeFactor
```

In differential modification mode, *exomePeak2* will export the differential modification peaks in formats of **BED** file and **CSV** table, the data will also be saved automatically under a folder named by `exomePeak2_output`.

Explanation for the additional table columns:

- ***ModLog2FC_control*** : the modification log2 fold enrichment in the control condition.
- ***ModLog2FC_treated*** : the modification log2 fold enrichment in the treatment condition.
- ***DiffModLog2FC*** : the log2 Fold Change of differential modification.

- **pvalue** : the p value of the differential modification.
- **padj** : the adjusted p value using BH approach.

3 Quantification and Statistical Analysis with Single Based Modification Annotation

exomePeak2 supports the modification quantification and differential modification analysis on single based modification annotation. The modification sites with single based resolution can provide a more accurate mapping of modification locations compared with the peaks called directly from the MeRIP-seq datasets.

Some of the datasets in epitranscriptomics have a single based resolution, e.x. Data generated by the m6A-CLIP-seq or m6A-miCLIP-seq techniques. exomePeak2 could provide a more accurate and consistent quantification and modification status inference for MeRIP-seq experiments using single based annotation.

exomePeak2 will automatically initiate the mode of single based modification quantification by providing a single based annotation file under the argument `mod_annot`.

The single based annotation information should be provided to the exomePeak2 function in the format of a `GRanges` object.

```
f2 = system.file("extdata", "mod_annot.rds", package="exomePeak2")

MOD_anno_GRANGE <- readRDS(f2)

exomePeak2(bam_ip = IP_BAM,
            bam_input = INPUT_BAM,
            gff_dir = GENE_anno_GTF,
            genome = "hg19",
            paired_end = FALSE,
            mod_annot = MOD_anno_GRANGE)

## class: SummarizedExomePeak
## dim: 171 7
## metadata(0):
## assays(2): ' GCsizeFactors
## rownames(171): mod_1 mod_2 ... control_83 control_84
## rowData names(2): GC_content feature_length
## colnames(7): IP1.bam IP2.bam ... Input2.bam Input3.bam
## colData names(3): design_IP design_Treatment sizeFactor
```

In this mode, exomePeak2 will export the analysis result also in formats of **BED** file and **CSV** table, while each row of the table corresponds to the sites of the annotation `GRanges`.

4 Peak Calling and Visualization in Multiple Steps

The exomePeak2 package can achieve peak calling and peak statistics calculation with multiple functions.

1. Check the bam files of MeRIP-seq data before peak calling.

The *exomePeak2* user's guide

```
MeRIP_Seq_Alignment <- scanMeripBAM(  
  bam_ip = IP_BAM,  
  bam_input = INPUT_BAM,  
  paired_end = FALSE  
)
```

For MeRIP-seq experiment with interactive design (contain control and treatment groups), use the following code.

```
MeRIP_Seq_Alignment <- scanMeripBAM(  
  bam_ip = IP_BAM,  
  bam_input = INPUT_BAM,  
  bam_treated_input = TREATED_INPUT_BAM,  
  bam_treated_ip = TREATED_IP_BAM,  
  paired_end = FALSE  
)
```

2. Conduct peak calling analysis on exons using the provided bam files.

```
SummarizedExomePeaks <- exomePeakCalling(merip_bams = MeRIP_Seq_Alignment,  
  gff_dir = GENE_ANNO_GTF,  
  genome = "hg19")
```

Alternatively, use the following code to quantify MeRIP-seq data on single based modification annotation.

```
SummarizedExomePeaks <- exomePeakCalling(merip_bams = MeRIP_Seq_Alignment,  
  gff_dir = GENE_ANNO_GTF,  
  genome = "hg19",  
  mod_annot = MOD_ANNO_GRANGE)
```

3. Estimate size factors that are required for GC content bias correction.

```
SummarizedExomePeaks <- normalizeGC(SummarizedExomePeaks)
```

4. Report the statistics of modification peaks using Generalized Linear Model (GLM).

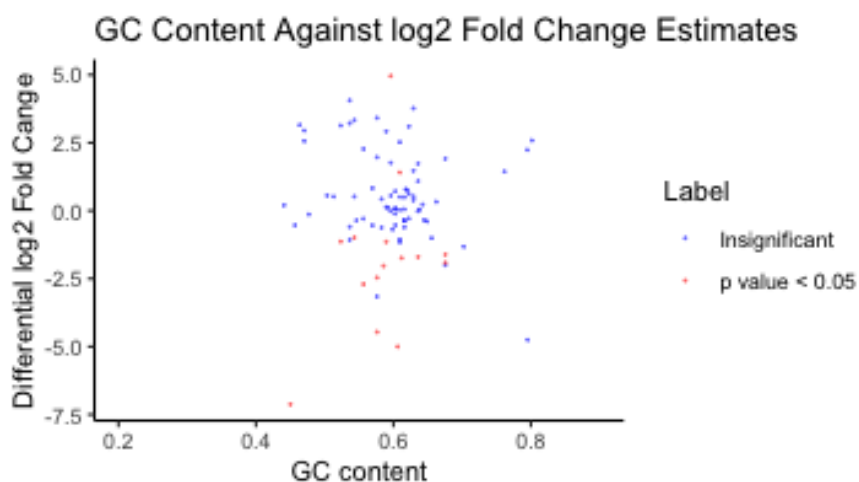
```
SummarizedExomePeaks <- glmM(SummarizedExomePeaks)
```

Alternatively, If the treated IP and input bam files are provided, `glmDM` function could be used to conduct differential modification analysis on modification Peaks with interactive GLM.

```
SummarizedExomePeaks <- glmDM(SummarizedExomePeaks)
```

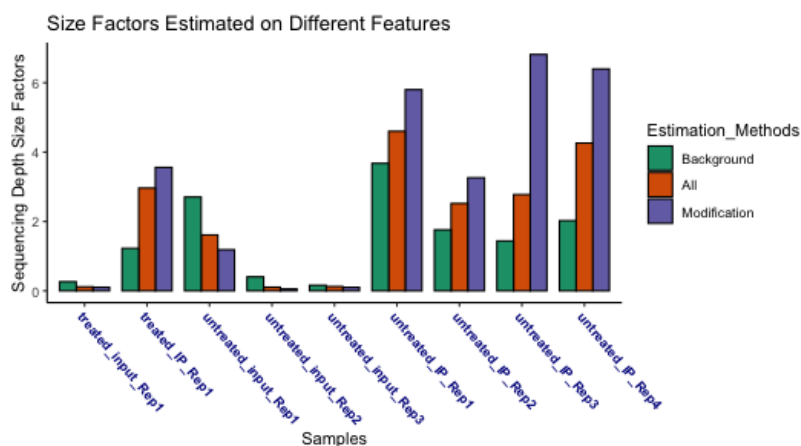
5. Generate the scatter plot between GC content and log2 Fold Change (LFC).

```
plotLfcGC(SummarizedExomePeaks)
```



6. Generate the bar plot for the sequencing depth size factors.

```
plotSizeFactors(SummarizedExomePeaks)
```



7. Export the modification peaks and the peak statistics with user decided format.

```
exportResults(SummarizedExomePeaks, format = "BED")
```

5 Contact

Please contact the maintainer of *exomePeak2* if you have encountered any problems:

ZhenWei : zhen.wei@xjtlu.edu.cn

Please visit the github page of *exomePeak2*:

<https://github.com/ZhenWei10/exomePeak2>

6 Session Info

```
sessionInfo()
## R version 3.5.3 (2019-03-11)
## Platform: x86_64-apple-darwin15.6.0 (64-bit)
## Running under: macOS Sierra 10.12.6
##
## Matrix products: default
## BLAS: /Library/Frameworks/R.framework/Versions/3.5/Resources/lib/libRblas.0.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/3.5/Resources/lib/libRlapack.dylib
##
## locale:
## [1] zh_CN.UTF-8/zh_CN.UTF-8/zh_CN.UTF-8/C/zh_CN.UTF-8/zh_CN.UTF-8
##
## attached base packages:
## [1] splines    parallel  stats4    stats      graphics  grDevices  utils
## [8] datasets  methods   base
##
## other attached packages:
## [1] BSgenome.Hsapiens.UCSC.hg19_1.4.0 BSgenome_1.50.0
## [3] rtracklayer_1.42.2                  Biostrings_2.50.2
## [5] XVector_0.22.0                      exomePeak2_0.9.9
## [7] knitr_1.23                          cqn_1.28.1
## [9] quantreg_5.51                       SparseM_1.77
## [11] preprocessCore_1.44.0               nor1mix_1.3-0
## [13] mclust_5.4.5                       SummarizedExperiment_1.12.0
## [15] DelayedArray_0.8.0                 BiocParallel_1.16.6
## [17] matrixStats_0.54.0                 Biobase_2.42.0
## [19] GenomicRanges_1.34.0               GenomeInfoDb_1.18.2
## [21] IRanges_2.16.0                     S4Vectors_0.20.1
## [23] BiocGenerics_0.28.0                 BiocStyle_2.10.0
##
## loaded via a namespace (and not attached):
## [1] colorspace_1.4-1                   htmlTable_1.13.1
## [3] base64enc_0.1-3                    rstudioapi_0.10
## [5] MatrixModels_0.4-1                 bit64_0.9-7
## [7] AnnotationDbi_1.44.0               apeglm_1.4.2
## [9] geneplotter_1.60.0                 zeallot_0.1.0
## [11] Formula_1.2-3                      Rsamtools_1.34.1
## [13] annotate_1.60.1                     cluster_2.1.0
## [15] BiocManager_1.30.4                 compiler_3.5.3
## [17] http_1.4.0                          backports_1.1.5
## [19] assertthat_0.2.1                   Matrix_1.2-17
## [21] lazyeval_0.2.2                     acepack_1.4.1
## [23] htmltools_0.3.6                    prettyunits_1.0.2
## [25] tools_3.5.3                         coda_0.19-3
## [27] gtable_0.3.0                       glue_1.3.1
## [29] GenomeInfoDbData_1.2.0             reshape2_1.4.3
## [31] dplyr_0.8.3                        Rcpp_1.0.2
## [33] bbmle_1.0.20                       vctrs_0.2.0
```

The *exomePeak2* user's guide

```
## [35] xfun_0.8                stringr_1.4.0
## [37] XML_3.98-1.20           zlibbioc_1.28.0
## [39] MASS_7.3-51.4           scales_1.0.0
## [41] hms_0.5.0               RMariaDB_1.0.6
## [43] RColorBrewer_1.1-2      yaml_2.2.0
## [45] memoise_1.1.0           gridExtra_2.3
## [47] ggplot2_3.2.1           emdbook_1.3.11
## [49] biomaRt_2.38.0          rpart_4.1-15
## [51] latticeExtra_0.6-28     stringi_1.4.3
## [53] RSQLite_2.1.2           genefilter_1.64.0
## [55] checkmate_1.9.4         GenomicFeatures_1.34.8
## [57] rlang_0.4.1             pkgconfig_2.0.3
## [59] bitops_1.0-6            evaluate_0.14
## [61] lattice_0.20-38         purrr_0.3.2
## [63] labeling_0.3            GenomicAlignments_1.18.1
## [65] htmlwidgets_1.3         bit_1.1-14
## [67] tidyselect_0.2.5        plyr_1.8.4
## [69] magrittr_1.5            bookdown_0.12
## [71] DESeq2_1.22.2           R6_2.4.0
## [73] Hmisc_4.2-0            DBI_1.0.0
## [75] pillar_1.4.2            foreign_0.8-71
## [77] survival_2.44-1.1      RCurl_1.95-4.12
## [79] nnet_7.3-12            tibble_2.1.3
## [81] crayon_1.3.4            rmarkdown_1.14
## [83] progress_1.2.2          locfit_1.5-9.1
## [85] grid_3.5.3             data.table_1.12.2
## [87] blob_1.2.0             digest_0.6.22
## [89] xtable_1.8-4           numDeriv_2016.8-1.1
## [91] munsell_0.5.0
```