

An Introduction to *Guitar* Package

Xiao Du

Modified: 26 April, 2019. Compiled: May 2, 2019

1 Quick Start with Guitar

This is a manual for Guitar package. The Guitar package is aimed for RNA landmark-guided transcriptomic analysis of RNA-related genomic features.

The Guitar package enables the comparison of multiple genomic features, which need to be stored in a name list. Please see the following example, which reads 1000 RNA m6A methylation sites into R for detection. Of course, in actual data analysis, features may come from multiple sets of resources.

```
library(Guitar)

## Loading required package: GenomicFeatures
## Loading required package: BiocGenerics
## Loading required package: parallel
##
## Attaching package: 'BiocGenerics'
## The following objects are masked from 'package:parallel':
##
##   clusterApply, clusterApplyLB, clusterCall,
##   clusterEvalQ, clusterExport, clusterMap,
##   parApply, parCapply, parLapply, parLapplyLB,
##   parRapply, parSapply, parSapplyLB
## The following objects are masked from 'package:stats':
##
##   IQR, mad, sd, var, xtabs
## The following objects are masked from 'package:base':
##
##   Filter, Find, Map, Position, Reduce,
##   anyDuplicated, append, as.data.frame, basename,
##   cbind, colnames, dirname, do.call, duplicated,
##   eval, evalq, get, grep, grepl, intersect,
##   is.unsorted, lapply, mapply, match, mget, order,
##   paste, pmax, pmax.int, pmin, pmin.int, rank,
##   rbind, rownames, sapply, setdiff, sort, table,
##   tapply, union, unique, unsplit, which, which.max,
##   which.min
## Loading required package: S4Vectors
## Loading required package: stats4
##
## Attaching package: 'S4Vectors'
## The following object is masked from 'package:base':
##
##   expand.grid
```

```

## Loading required package: IRanges
## Loading required package: GenomeInfoDb
## Loading required package: GenomicRanges
## Loading required package: AnnotationDbi
## Loading required package: Biobase
## Welcome to Bioconductor
##
##   Vignettes contain introductory material; view
##   with 'browseVignettes()'. To cite Bioconductor,
##   see 'citation("Biobase")', and for packages
##   'citation("pkgname")'.
## Loading required package: rtracklayer
## Loading required package: magrittr
## Loading required package: ggplot2
## Registered S3 methods overwritten by 'ggplot2':
##   method      from
## [.quosures    rlang
## c.quosures     rlang
## print.quosures rlang
## Loading required package: dplyr
##
## Attaching package: 'dplyr'
## The following object is masked from 'package:AnnotationDbi':
##
##   select
## The following object is masked from 'package:Biobase':
##
##   combine
## The following objects are masked from 'package:GenomicRanges':
##
##   intersect, setdiff, union
## The following object is masked from 'package:GenomeInfoDb':
##
##   intersect
## The following objects are masked from 'package:IRanges':
##
##   collapse, desc, intersect, setdiff, slice, union
## The following objects are masked from 'package:S4Vectors':
##
##   first, intersect, rename, setdiff, setequal,
##   union
## The following objects are masked from 'package:BiocGenerics':
##
##   combine, intersect, setdiff, union
## The following objects are masked from 'package:stats':
##
##   filter, lag
## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union
##
## Attaching package: 'Guitar'
## The following object is masked from 'package:BiocGenerics':
##

```

```
##      normalize

# genomic features imported into named list
stBedFiles <- list(system.file("extdata", "m6A_mm10_exomePeak_1000peaks_bed12.bed",
                             package="Guitar"))
```

With the following script, we may generate the transcriptomic distribution of genomic features to be tested, and the result will be automatically saved into a PDF file under the working directory with prefix "example". With the `GuitarPlot` function, the gene annotation can be downloaded from internet automatically with a genome assembly number provided; however, this feature requires working internet and might take a longer time. The toy Guitar coordinates generated internally should never be re-used in other real data analysis.

```
count <- GuitarPlot(txGenomeVer = "mm10",
                   stBedFiles = stBedFiles,
                   miscOutFilePrefix = NA)
```

In a more efficient protocol, in order to re-use the gene annotation and *Guitar coordinates*, you will have to build Guitar Coordinates from a *txdb* object in a separate step. The *transcriptDb* contains the gene annotation information and can be obtained in a number of ways, .e.g, download the complete gene annotation of species from UCSC automatically, which might takes a few minutes. In the following analysis, we load the *Txdb* object from a toy dataset provided with the Guitar package. Please note that this is only a very small part of the complete hg19 transcriptome, and the *Txdb* object provided with *Guitar* package should not be used in real data analysis. With a *Txdb* object that contains gene annotation information, we in the next build *Guitar coordinates*, which is essentially a bridge connects the transcriptomic landmarks and genomic coordinates.

```
txdb_file <- system.file("extdata", "mm10_toy.sqlite",
                        package="Guitar")
txdb <- loadDb(txdb_file)
guitarTxdb <- makeGuitarTxdb(txdb = txdb, txPrimaryOnly = FALSE)

## [1] "There are 2946 transcripts of 2946 genes in the genome."
## [1] "total 2946 transcripts extracted ..."
## [1] "total 2719 transcripts left after ambiguity filter ..."
## [1] "total 2719 transcripts left after check chromosome validity ..."
## [1] "total 1342 mRNAs left after component length filter ..."
## [1] "total 307 ncRNAs left after ncRNA length filter ..."
## [1] "generate components for all tx"
## [1] "generate components for mRNA"
## [1] "generate components for lncRNA"
## [1] "generate chiped transcriptome"
## [1] "generate coverage checking ranges for tx"
## [1] "generate coverage checking ranges for mrna"
## [1] "generate coverage checking ranges for ncrna"

# Or use gff. file to generate guitarTxdb
# Or use getTxdb() to download TxDb from internet:
# txdb <- getTxdb(txGenomeVer="hg19")
# guitarTxdb <- makeGuitarTxdb(txdb)
```

You may now generate the Guitar plot from the named list of genome-based features.

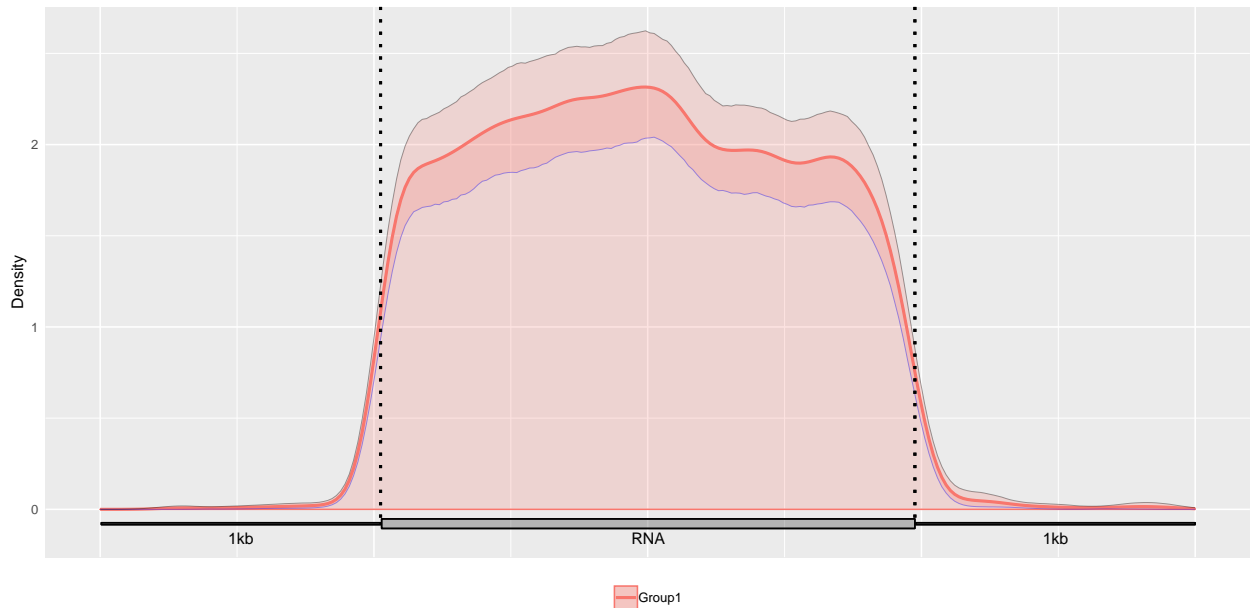
```
GuitarPlot(txTxdb = txdb,
            stBedFiles = stBedFiles,
            miscOutFilePrefix = "example")

## [1] "20190502224421"
## [1] "There are 2946 transcripts of 2946 genes in the genome."
## [1] "total 2946 transcripts extracted ..."
## [1] "total 2719 transcripts left after ambiguity filter ..."
## [1] "total 2719 transcripts left after check chromosome validity ..."
## [1] "total 1342 mRNAs left after component length filter ..."
## [1] "total 307 ncRNAs left after ncRNA length filter ..."
## [1] "generate components for all tx"
## [1] "generate components for mRNA"
## [1] "generate components for lncRNA"
## [1] "generate chiped transcriptome"
## [1] "generate coverage checking ranges for tx"
## [1] "generate coverage checking ranges for mrna"
## [1] "generate coverage checking ranges for ncrna"
## [1] "20190502224435"
## [1] "import BED file /private/tmp/RtmpC6NAv7/Rinst121819830c0e/Guitar/extdata/m6A_mm10_exomePeak_1000
## [1] "sample 10 points for Group1"
## [1] "start figure plotting for tx ..."
## [1] "start figure plotting for mrna ..."
## [1] "start figure plotting for ncrna ..."
```

Alternatively, you may also optionally include the promoter DNA region and tail DNA region on the 5' and 3' side of a transcript in the plot with parameter `headOrtail = TRUE`.

```
GuitarPlot(txTxdb = txdb,
            stBedFiles = stBedFiles,
            headOrtail = TRUE)

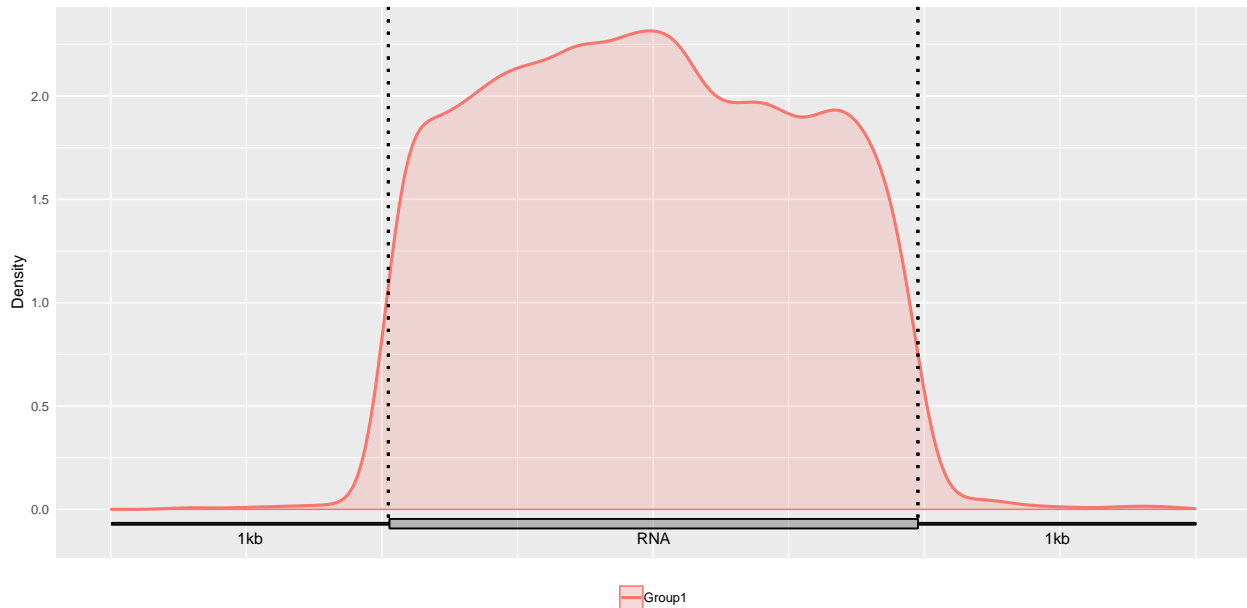
## [1] "20190502224534"
## [1] "There are 2946 transcripts of 2946 genes in the genome."
## [1] "total 2946 transcripts extracted ..."
## [1] "total 2719 transcripts left after ambiguity filter ..."
## [1] "total 2719 transcripts left after check chromosome validity ..."
## [1] "total 1342 mRNAs left after component length filter ..."
## [1] "total 307 ncRNAs left after ncRNA length filter ..."
## [1] "generate components for all tx"
## [1] "generate components for mRNA"
## [1] "generate components for lncRNA"
## [1] "generate chiped transcriptome"
## [1] "generate coverage checking ranges for tx"
## [1] "generate coverage checking ranges for mrna"
## [1] "generate coverage checking ranges for ncrna"
## [1] "20190502224550"
## [1] "import BED file /private/tmp/RtmpC6NAv7/Rinst121819830c0e/Guitar/extdata/m6A_mm10_exomePeak_1000
## [1] "sample 10 points for Group1"
## [1] "start figure plotting for tx ..."
```



Alternatively, you may also optionally include the Confidence Interval for guitar plot with parameter `enableCI = FALSE`.

```
GuitarPlot(txTxdb = txdb,
           stBedFiles = stBedFiles,
           headOrtail = TRUE,
           enableCI = FALSE)

## [1] "20190502224623"
## [1] "There are 2946 transcripts of 2946 genes in the genome."
## [1] "total 2946 transcripts extracted ..."
## [1] "total 2719 transcripts left after ambiguity filter ..."
## [1] "total 2719 transcripts left after check chromosome validity ..."
## [1] "total 1342 mRNAs left after component length filter ..."
## [1] "total 307 ncRNAs left after ncRNA length filter ..."
## [1] "generate components for all tx"
## [1] "generate components for mRNA"
## [1] "generate components for lncRNA"
## [1] "generate chiped transcriptome"
## [1] "generate coverage checking ranges for tx"
## [1] "generate coverage checking ranges for mrna"
## [1] "generate coverage checking ranges for ncrna"
## [1] "20190502224638"
## [1] "import BED file /private/tmp/RtmpC6Nav7/Rinst121819830c0e/Guitar/extdata/m6A_mm10_exomePeak_1000"
## [1] "sample 10 points for Group1"
## [1] "start figure plotting for tx ..."
```



2 Supported Data Format

Besides BED file, Guitar package also supports GRangesList and GRanges data structures. Please see the following examples.

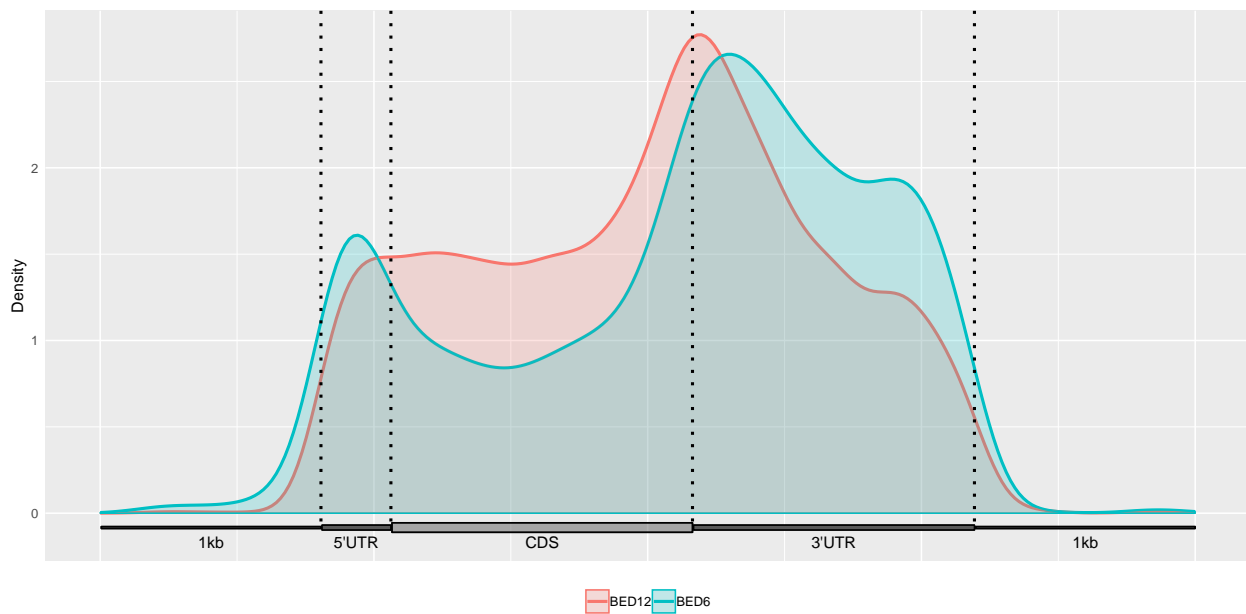
```
# import different data formats into a named list object.
# These genomic features are using mm10 genome assembly
stBedFiles <- list(system.file("extdata", "m6A_mm10_exomePeak_1000peaks_bed12.bed",
                             package="Guitar"),
                  system.file("extdata", "m6A_mm10_exomePeak_1000peaks_bed6.bed",
                             package="Guitar"))

# Build Guitar Coordinates
txdb_file <- system.file("extdata", "mm10_toy.sqlite",
                        package="Guitar")
txdb <- loadDb(txdb_file)

# Guitar Plot
GuitarPlot(txTxdb = txdb,
           stBedFiles = stBedFiles,
           headOrtail = TRUE,
           enableCI = FALSE,
           mapFilterTranscript = TRUE,
           pltTxType = c("mrna"),
           stGroupName = c("BED12", "BED6"))

## [1] "20190502224642"
## [1] "There are 2946 transcripts of 2946 genes in the genome."
## [1] "total 2946 transcripts extracted ..."
## [1] "total 2719 transcripts left after ambiguity filter ..."
## [1] "total 2719 transcripts left after check chromosome validity ..."
## [1] "total 1342 mRNAs left after component length filter ..."
## [1] "total 307 ncRNAs left after ncRNA length filter ..."
```

```
## [1] "generate components for mRNA"
## [1] "generate chiped transcriptome"
## [1] "generate coverage checking ranges for mrna"
## [1] "20190502224658"
## [1] "import BED file /private/tmp/RtmpC6NAv7/Rinst121819830c0e/Guitar/extdata/m6A_mm10_exomePeak_1000"
## [1] "import BED file /private/tmp/RtmpC6NAv7/Rinst121819830c0e/Guitar/extdata/m6A_mm10_exomePeak_1000"
## [1] "sample 10 points for BED12"
## [1] "sample 10 points for BED6"
## [1] "start figure plotting for mrna ..."
```



3 Processing of sampling sites information

We can select parameters for site sampling.

```
stGRangeLists = vector("list", length(stBedFiles))
sitesPoints <- list()
for (i in seq_len(length(stBedFiles))) {
  stGRangeLists[[i]] <- blocks(import(stBedFiles[[i]]))
}
for (i in seq_len(length(stGRangeLists))) {
  sitesPoints[[i]] <- samplePoints(stGRangeLists[i],
    stSampleNum = 10,
    stAmbiguity = 5,
    pltTxType = c("mrna"),
    stSampleModle = "Equidistance",
    mapFilterTranscript = FALSE,
    guitarTxdb = guitarTxdb)
}
```

4 Guitar Coordinates - Transcriptomic Landmarks Projected on Genome

The `guitarTxdb` object contains the genome-projected transcriptome coordinates, which can be valuable for evaluating transcriptomic information related applications, such as checking the quality of MeRIP-Seq data. The `Guitar` coordinates are essentially the genomic projection of standardized transcript-based coordinates, making a viable bridge between the landmarks on transcript and genome-based coordinates.

It is based on the `txdb` object input, extracts the transcript information in `txdb`, selects the transcripts that match the parameters according to the component parameters set by the user, and saves according to the transcript type (`tx`, `mrna`, `ncrna`).

```
guitarTxdb <- makeGuitarTxdb(txdb = txdb,
                             txAmbiguity = 5,
                             txMrnaComponentProp = c(0.1,0.15,0.6,0.05,0.1),
                             txLncrnaComponentProp = c(0.2,0.6,0.2),
                             pltTxType = c("tx","mrna","ncrna"),
                             txPrimaryOnly = FALSE)

## [1] "There are 2946 transcripts of 2946 genes in the genome."
## [1] "total 2946 transcripts extracted ..."
## [1] "total 2719 transcripts left after ambiguity filter ..."
## [1] "total 2719 transcripts left after check chromosome validity ..."
## [1] "total 1342 mRNAs left after component length filter ..."
## [1] "total 307 ncRNAs left after ncRNA length filter ..."
## [1] "generate components for all tx"
## [1] "generate components for mRNA"
## [1] "generate components for lncRNA"
## [1] "generate chiped transcriptome"
## [1] "generate coverage checking ranges for tx"
## [1] "generate coverage checking ranges for mrna"
## [1] "generate coverage checking ranges for ncrna"
```

5 Check the Overlapping between Different Components

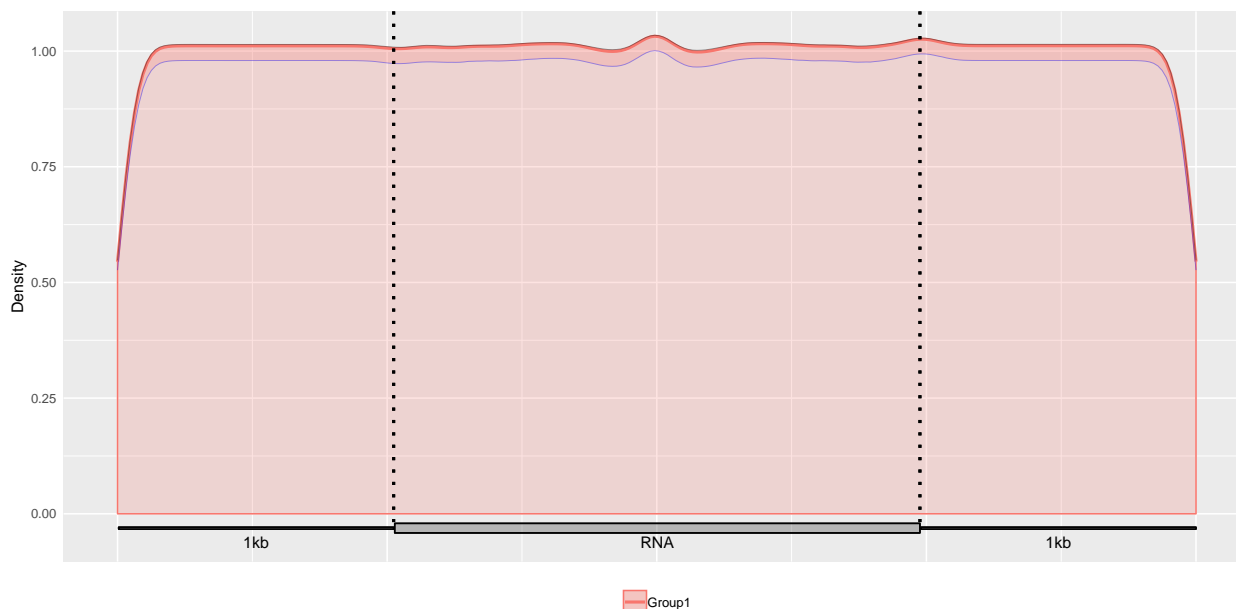
We can also check the distribution of the `Guitar` coordinates built.

```
gcl <- list(guitarTxdb$tx$tx)
GuitarPlot(txTxdb = txdb,
            stGRangeLists = gcl,
            stSampleNum = 200,
            enableCI = TRUE,
            pltTxType = c("tx"),
            txPrimaryOnly = FALSE
)

## [1] "20190502224719"
## [1] "There are 2946 transcripts of 2946 genes in the genome."
## [1] "total 2946 transcripts extracted ..."
## [1] "total 2719 transcripts left after ambiguity filter ..."
## [1] "total 2719 transcripts left after check chromosome validity ..."
## [1] "total 1342 mRNAs left after component length filter ..."
## [1] "total 307 ncRNAs left after ncRNA length filter ..."
## [1] "generate components for all tx"
## [1] "generate chiped transcriptome"
```



```
## [1] "generate coverage checking ranges for tx"
## [1] "20190502224734"
## [1] "sample 200 points for Group1"
## [1] "start figure plotting for tx ..."
```

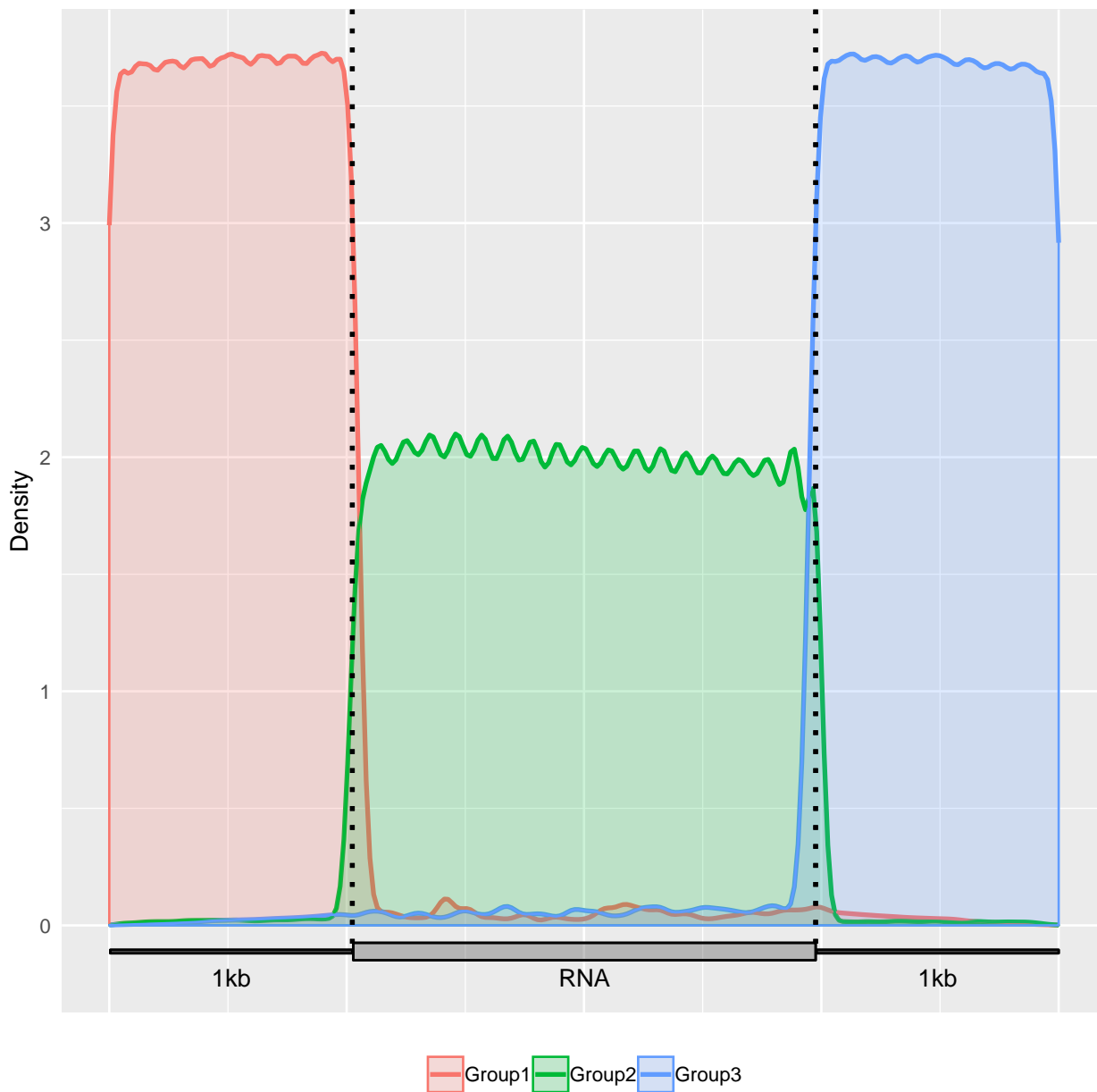


Alternatively, we can extract the RNA components, check the distribution of tx components in the transcriptome

```
GuitarCoords <- guitarTxdb$tx$txComponentGRange
type <- paste(mcols(GuitarCoords)$componentType, mcols(GuitarCoords)$txType)
key <- unique(type)
landmark <- list(1,2,3,4,5,6,7,8,9,10,11)
names(landmark) <- key
for (i in 1:length(key)) {
  landmark[[i]] <- GuitarCoords[type==key[i]]
}
GuitarPlot(txTxdb = txdb ,
            stGRangeLists = landmark[1:3],
            pltTxType = c("tx"),
            enableCI = FALSE
)

## [1] "20190502230209"
## [1] "There are 2946 transcripts of 2946 genes in the genome."
## [1] "total 2946 transcripts extracted ..."
## [1] "total 2719 transcripts left after ambiguity filter ..."
## [1] "total 2719 transcripts left after check chromosome validity ..."
## [1] "total 1342 mRNAs left after component length filter ..."
## [1] "total 307 ncRNAs left after ncRNA length filter ..."
## [1] "generate components for all tx"
## [1] "generate chiped transcriptome"
## [1] "generate coverage checking ranges for tx"
## [1] "20190502230224"
## [1] "sample 10 points for Group1"
```

```
## [1] "sample 10 points for Group2"
## [1] "sample 10 points for Group3"
## [1] "start figure plotting for tx ..."
```



Check the distribution of mRNA components in the transcriptome

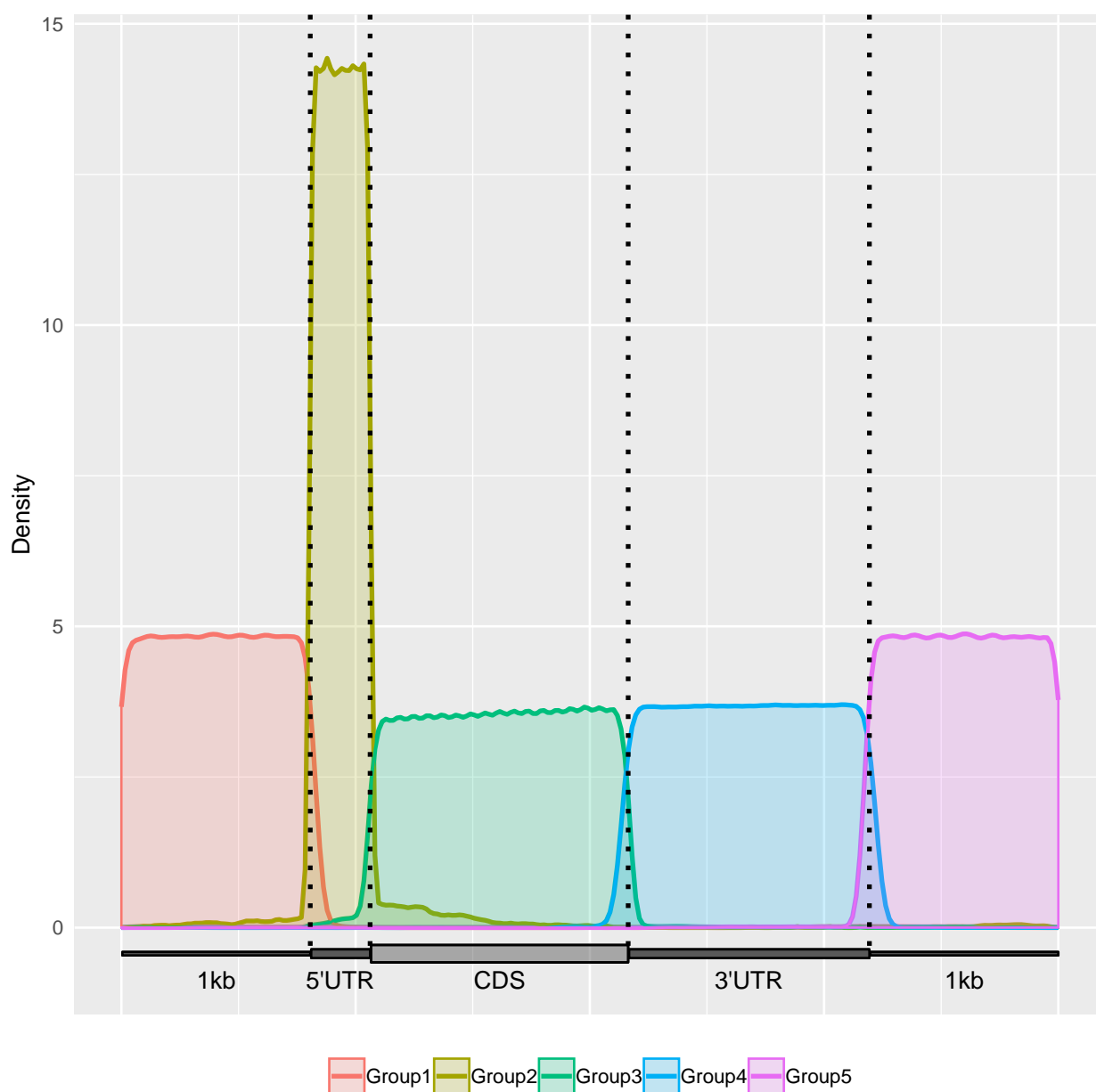
```
GuitarPlot(txTxdb = txdb ,
            stGRangeLists = landmark[4:8],
            pltTxType = c("mrna"),
            enableCI = FALSE
)

## [1] "20190502230304"
## [1] "There are 2946 transcripts of 2946 genes in the genome."
## [1] "total 2946 transcripts extracted ..."
```

```

## [1] "total 2719 transcripts left after ambiguity filter ..."
## [1] "total 2719 transcripts left after check chromosome validity ..."
## [1] "total 1342 mRNAs left after component length filter ..."
## [1] "total 307 ncRNAs left after ncRNA length filter ..."
## [1] "generate components for mRNA"
## [1] "generate chiped transcriptome"
## [1] "generate coverage checking ranges for mrna"
## [1] "20190502230318"
## [1] "sample 10 points for Group1"
## [1] "sample 10 points for Group2"
## [1] "sample 10 points for Group3"
## [1] "sample 10 points for Group4"
## [1] "sample 10 points for Group5"
## [1] "start figure plotting for mrna ..."

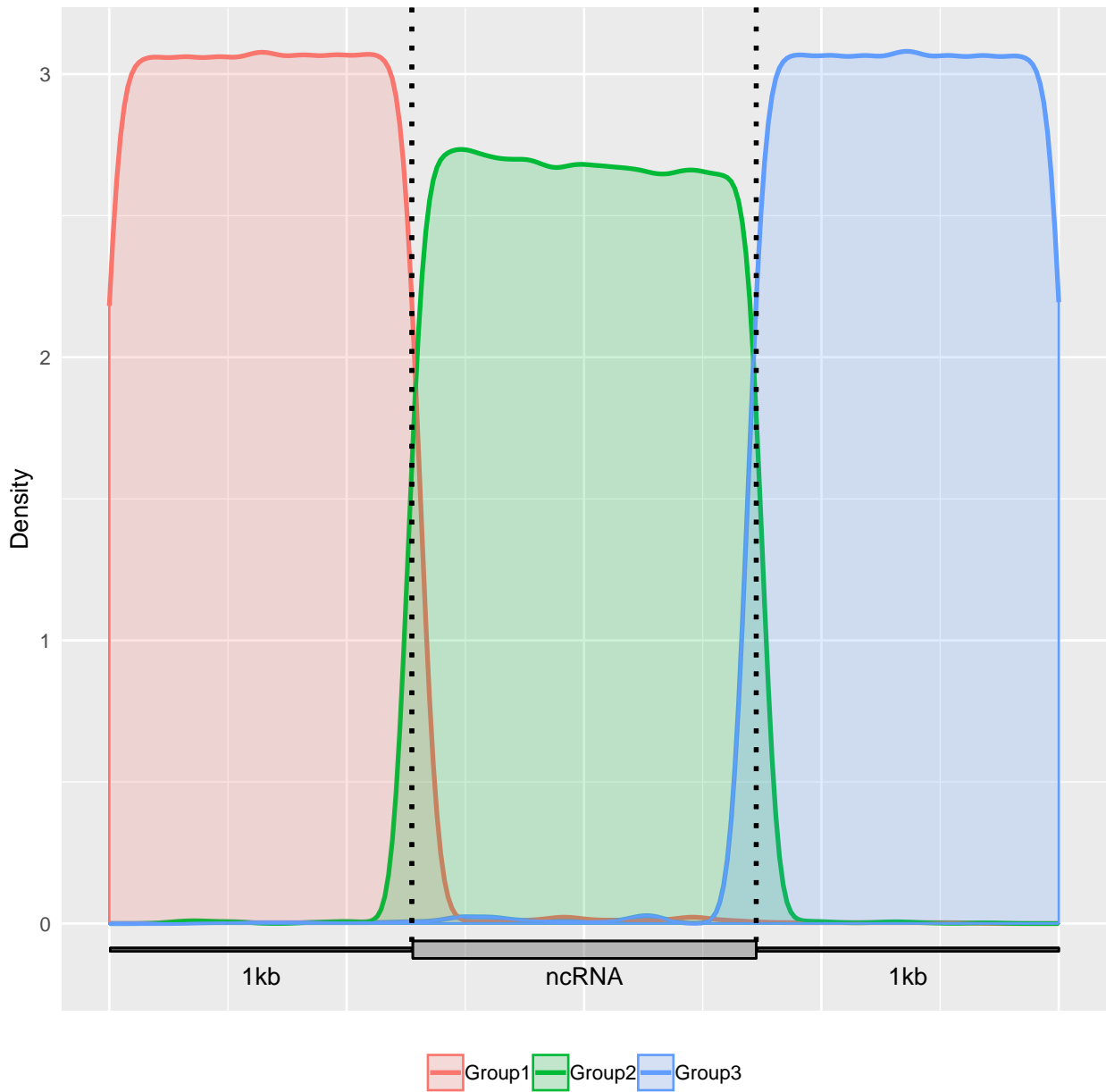
```



Check the distribution of lncRNA components in the transcriptome

```
GuitarPlot(txTxdb = txdb ,
           stGRangeLists = landmark[9:11],
           pltTxType = c("ncrna"),
           enableCI = FALSE
)

## [1] "20190502230347"
## [1] "There are 2946 transcripts of 2946 genes in the genome."
## [1] "total 2946 transcripts extracted ..."
## [1] "total 2719 transcripts left after ambiguity filter ..."
## [1] "total 2719 transcripts left after check chromosome validity ..."
## [1] "total 1342 mRNAs left after component length filter ..."
## [1] "total 307 ncRNAs left after ncRNA length filter ..."
## [1] "generate components for lncRNA"
## [1] "generate chiped transcriptome"
## [1] "generate coverage checking ranges for ncrna"
## [1] "20190502230410"
## [1] "sample 10 points for Group1"
## [1] "sample 10 points for Group2"
## [1] "sample 10 points for Group3"
## [1] "start figure plotting for ncrna ..."
```



6 Session Information

```
sessionInfo()

## R version 3.6.0 (2019-04-26)
## Platform: x86_64-apple-darwin15.6.0 (64-bit)
## Running under: OS X El Capitan 10.11.6
##
## Matrix products: default
## BLAS: /Library/Frameworks/R.framework/Versions/3.6/Resources/lib/libRblas.0.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/3.6/Resources/lib/libRlapack.dylib
##
```

```

## locale:
## [1] C/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## attached base packages:
## [1] stats4      parallel  stats      graphics  grDevices
## [6] utils       datasets  methods    base
##
## other attached packages:
## [1] Guitar_2.0.0          dplyr_0.8.0.1
## [3] ggplot2_3.1.1         magrittr_1.5
## [5] rtracklayer_1.43.3    GenomicFeatures_1.36.0
## [7] AnnotationDbi_1.46.0  Biobase_2.44.0
## [9] GenomicRanges_1.36.0 GenomeInfoDb_1.20.0
## [11] IRanges_2.18.0        S4Vectors_0.22.0
## [13] BiocGenerics_0.30.0   knitr_1.22
##
## loaded via a namespace (and not attached):
## [1] Rcpp_1.0.1            lattice_0.20-38
## [3] prettyunits_1.0.2     Rsamtools_2.0.0
## [5] Biostrings_2.52.0     assertthat_0.2.1
## [7] digest_0.6.18         R6_2.4.0
## [9] plyr_1.8.4            RSQLite_2.1.1
## [11] evaluate_0.13         httr_1.4.0
## [13] highr_0.8             pillar_1.3.1
## [15] zlibbioc_1.30.0       rlang_0.3.4
## [17] progress_1.2.0        lazyeval_0.2.2
## [19] blob_1.1.1            Matrix_1.2-17
## [21] labeling_0.3          BiocParallel_1.18.0
## [23] stringr_1.4.0         RCurl_1.95-4.12
## [25] bit_1.1-14            biomaRt_2.40.0
## [27] munsell_0.5.0         DelayedArray_0.10.0
## [29] compiler_3.6.0        xfun_0.6
## [31] pkgconfig_2.0.2       tidyselect_0.2.5
## [33] SummarizedExperiment_1.14.0 tibble_2.1.1
## [35] GenomeInfoDbData_1.2.1 matrixStats_0.54.0
## [37] XML_3.98-1.19         crayon_1.3.4
## [39] withr_2.1.2           GenomicAlignments_1.20.0
## [41] bitops_1.0-6          grid_3.6.0
## [43] gtable_0.3.0          DBI_1.0.0
## [45] scales_1.0.0          stringi_1.4.3
## [47] XVector_0.24.0        tools_3.6.0
## [49] bit64_0.9-7           glue_1.3.1
## [51] purrr_0.3.2           hms_0.4.2
## [53] colorspace_1.4-1      memoise_1.1.0

```