

Assign mir181 binding sites to a specific transcript

Melina Klostermann

20 March, 2024

Contents

1	Libraries and settings	1
2	What was done?	1
3	Assign each binding site to a specific transcript	2
4	Get transcript coordinates	3
5	Session info	5

1 Libraries and settings

```
# -----  
# libraries  
# -----  
library(tidyverse)  
library(GenomicRanges)  
library(GenomicFeatures)  
  
# -----  
# settings  
# -----  
  
here <- here::here()  
  
out <- paste0(here, "/Figure4/03_assign_transcripts/")
```

2 What was done?

- Each binding site is assigned to a specific transcript isoform
- The following criteria are used to decide on the isoform:
 - 1) Previously assigned region (see binding site definitions), if more then one transcript is possible:
 - 2) if possible protein coding transcripts, if still more then one transcript is possible:
 - 3) best transcript support level, if still more then one transcript is possible:
 - 4) longest transcript

- Then the mir181 binding sites are mapped to their respective transcript coordinates.
- The transcript annotations are later used for motif discovery and structure predictions.

```
#-----
# Files
#-----
anno <- readRDS(paste0(here, "/Supporting_scripts/annotation_preprocessing/annotation.rds"))
anno$gene_id <- sub("\\..*", "", anno$gene_id)
anno$transcript_id <- sub("\\..*", "", anno$transcript_id)

mir181_bs <- readRDS(paste0(here, "/Figure4/01_MRE_bound_gene_and_bound_region/mir181_bs_afterFigure4B.rds"))
```

3 Assign each binding site to a specific transcript

```
# get appris transcripts (when there are multiple take the longest)
transcripts <- anno[anno$type=="transcript"] %>% as.data.frame(.)
#transcripts$appris <- grepl(transcripts$tag, pattern= "appris_principal_1")

# get regional annotations
three_utr <- anno[anno$type=="three_prime_UTR"] %>% as.data.frame(.) %>%
  dplyr::select(seqnames, start, end, width, strand, transcript_id, gene_id) %>%
  left_join(transcripts %>% dplyr::select(., width, transcript_id, transcript_type, transcript_support_level) %>%
    makeGRangesFromDataFrame(., keep.extra.columns = T))

cds <- anno[anno$type=="CDS"] %>% as.data.frame(.) %>%
  dplyr::select(seqnames, start, end, width, strand, transcript_id, gene_id) %>%
  left_join(transcripts %>% dplyr::select(., width, transcript_id, transcript_type, transcript_support_level) %>%
    makeGRangesFromDataFrame(., keep.extra.columns = T))

five_utr <- anno[anno$type=="five_prime_UTR"] %>% as.data.frame(.) %>%
  dplyr::select(seqnames, start, end, width, strand, transcript_id, gene_id) %>%
  left_join(transcripts %>% dplyr::select(., width, transcript_id, transcript_type, transcript_support_level) %>%
    makeGRangesFromDataFrame(., keep.extra.columns = T))

# select best transcript per bs
# 3'utr
bs_utr3 <- mir181_bs %>% subset(., region == "utr3") %>%
  makeGRangesFromDataFrame(., keep.extra.columns = T)
NROW(bs_utr3)

## [1] 5607

i_utr3 <- findOverlaps(bs_utr3, three_utr, type = "any")
t <- as.data.frame(three_utr[subjectHits(i_utr3)]) %>% dplyr::select(width_tx, transcript_type, transcript_support_level)
t$transcript_type <- factor(t$transcript_type, levels = c("protein_coding"))

bs_utr3 <- bs_utr3[queryHits(i_utr3)] %>%
  as.tibble()

bs_utr3 <- cbind(bs_utr3, t) %>%
  group_by(mir181BS_ID) %>%
  arrange(transcript_type, transcript_support_level, width_tx, .by_group = T) %>%
```

```

dplyr::slice(1)
nrow(bs_utr3)

## [1] 5607

# cds
bs_cds <- mir181_bs %>% subset(., region == "cds") %>%
  makeGRangesFromDataFrame(., keep.extra.columns = T)
NROW(bs_cds)

## [1] 4284

i_cds <- findOverlaps(bs_cds, cds, type = "any")
t <- as.data.frame(cds[subjectHits(i_cds)]) %>% dplyr::select(width_tx, transcript_type, transcript_support_level)

bs_cds <- bs_cds[queryHits(i_cds)] %>%
  as.tibble()

bs_cds <- cbind(bs_cds, t) %>%
  group_by(mir181BS_ID) %>%
  arrange(transcript_type, transcript_support_level, width_tx, .by_group = T) %>%
  dplyr::slice(1)
nrow(bs_cds)

## [1] 4284

# 5'utr
bs_utr5 <- mir181_bs %>% subset(., region == "utr5") %>%
  makeGRangesFromDataFrame(., keep.extra.columns = T)
NROW(bs_utr5)

## [1] 582

i_utr5 <- findOverlaps(bs_utr5, five_utr, type = "any")
t <- as.data.frame(five_utr[subjectHits(i_utr5)]) %>% dplyr::select(width_tx, transcript_type, transcript_support_level)

bs_utr5 <- bs_utr5[queryHits(i_utr5)] %>%
  as.tibble()

bs_utr5 <- cbind(bs_utr5, t) %>%
  group_by(mir181BS_ID) %>%
  arrange(transcript_type, transcript_support_level, width_tx, .by_group = T) %>%
  dplyr::slice(1)
nrow(bs_utr5)

## [1] 582

mir181_bs <- rbind(bs_utr3, bs_cds, bs_utr5)

```

4 Get transcript coordinates

```

#####
# BS sequence considering mature transcripts
#####
# prepare a txdb of expressed transcripts
anno_transcripts_exons <- anno[anno$type != "gene"]

```

```

anno_transcripts_exons$transcript_id <- sub("\\..*", "", anno_transcripts_exons$transcript_id)
anno_transcripts_GR_list <- anno_transcripts_exons %>%
  splitAsList(., f = .$transcript_id) %>%
  GRangesList(.)

txdb <- makeTxDbFromGRanges(unlist(anno_transcripts_GR_list))

# prepare a transcript mapper (contains transcript ids and names together with genomic positions of transcripts)
transcripts_txdb_mapper <- transcripts(txdb)

# get transcript-relative coordinates of BS
mir181_bs <- makeGRangesFromDataFrame(mir181_bs, keep.extra.columns = T)
mir181_bs_tx <- mapToTranscripts(mir181_bs, txdb, extractor.fun = GenomicFeatures::exonsBy, ignore.strand = T)
# Mapped position is computed by counting from the transcription start site (TSS) and is not affected by strand

# read metadata
gen_pos <- mir181_bs[mir181_bs_tx$xHits] %>%
  as.data.frame() %>%
  dplyr::select(start, end, strand, seqnames) %>%
  rename( start = "genomic_start" , end = "genomic_end", strand = "genomic_strand", seqnames = "genomic_seqnames")

elementMetadata(mir181_bs_tx) <- c(elementMetadata(mir181_bs_tx), elementMetadata(mir181_bs[mir181_bs_tx$transcriptsHits]))

# change the seqnames to the transcript names
names(mir181_bs_tx) <- 1: NROW(mir181_bs_tx)
mir181_bs_tx <- as.data.frame(mir181_bs_tx)
mir181_bs_tx$seqnames <- transcripts_txdb_mapper$tx_name[mir181_bs_tx$transcriptsHits]

mir181_bs_tx <- mir181_bs_tx %>% subset(seqnames == transcript_id)

```

- Number of mirBS: 10473
- Number of mirBS on transcripts (without intron): 10244
- Number of enriched mirBS: 4658
- Number of enriched mirBS on transcripts: 4519

Comment: we use findOverlaps with the center position of the binding site to assign the regions in the binding site definition scripts. The mapToTranscripts will only keep binding sites, that are completely inside the transcript. This is why we loose a few binding sites, when mapping to transcripts here.

```

saveRDS(mir181_bs_tx, paste0(out,"mir181_bs_on_transcripts.rds"))
saveDb(txdb, paste0(out,"transcript_annotation.db"))

```

```

## TxDb object:
## # Db type: TxDb
## # Supporting package: GenomicFeatures
## # Genome: NA
## # Nb of transcripts: 142351
## # Db created by: GenomicFeatures package from Bioconductor
## # Creation time: 2024-03-20 10:35:35 +0100 (Wed, 20 Mar 2024)
## # GenomicFeatures version at creation time: 1.54.3
## # RSQLite version at creation time: 2.3.5
## # DBSCHEMAVERSION: 1.2

```

5 Session info

```
sessionInfo()
```

```
## R version 4.3.2 (2023-10-31)
## Platform: x86_64-apple-darwin20 (64-bit)
## Running under: macOS Ventura 13.1
##
## Matrix products: default
## BLAS: /Library/Frameworks/R.framework/Versions/4.3-x86_64/Resources/lib/libRblas.0.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/4.3-x86_64/Resources/lib/libRlapack.dylib; LAPACK
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## time zone: Europe/Berlin
## tzcode source: internal
##
## attached base packages:
## [1] stats4 stats graphics grDevices utils datasets methods
## [8] base
##
## other attached packages:
## [1] GenomicFeatures_1.54.3 AnnotationDbi_1.64.1 Biobase_2.62.0
## [4] GenomicRanges_1.54.1 GenomeInfoDb_1.38.6 IRanges_2.36.0
## [7] S4Vectors_0.40.2 BiocGenerics_0.48.1 lubridate_1.9.3
## [10] forcats_1.0.0 stringr_1.5.1 dplyr_1.1.4
## [13] purrr_1.0.2 readr_2.1.5 tidyr_1.3.1
## [16] tibble_3.2.1 ggplot2_3.4.4 tidyverse_2.0.0
## [19] knitr_1.45
##
## loaded via a namespace (and not attached):
## [1] tidyselect_1.2.0 blob_1.2.4
## [3] filelock_1.0.3 Biostrings_2.70.2
## [5] bitops_1.0-7 fastmap_1.1.1
## [7] RCurl_1.98-1.14 BiocFileCache_2.10.1
## [9] GenomicAlignments_1.38.2 XML_3.99-0.16.1
## [11] digest_0.6.34 timechange_0.3.0
## [13] lifecycle_1.0.4 KEGGREST_1.42.0
## [15] RSQLite_2.3.5 magrittr_2.0.3
## [17] compiler_4.3.2 rlang_1.1.3
## [19] progress_1.2.3 tools_4.3.2
## [21] utf8_1.2.4 yaml_2.3.8
## [23] rtracklayer_1.62.0 S4Arrays_1.2.0
## [25] prettyunits_1.2.0 bit_4.0.5
## [27] curl_5.2.0 here_1.0.1
## [29] DelayedArray_0.28.0 xml2_1.3.6
## [31] abind_1.4-5 BiocParallel_1.36.0
## [33] withr_3.0.0 grid_4.3.2
## [35] fansi_1.0.6 colorspace_2.1-0
## [37] scales_1.3.0 SummarizedExperiment_1.32.0
## [39] biomaRt_2.58.2 cli_3.6.2
## [41] rmarkdown_2.25 crayon_1.5.2
## [43] generics_0.1.3 rstudioapi_0.15.0
```

## [45] httr_1.4.7	tzdb_0.4.0
## [47] rjson_0.2.21	DBI_1.2.1
## [49] cachem_1.0.8	zlibbioc_1.48.0
## [51] parallel_4.3.2	XVector_0.42.0
## [53] restfulr_0.0.15	matrixStats_1.2.0
## [55] vctrs_0.6.5	Matrix_1.6-5
## [57] hms_1.1.3	bit64_4.0.5
## [59] glue_1.7.0	codetools_0.2-19
## [61] stringi_1.8.3	gtable_0.3.4
## [63] BiocIO_1.12.0	munsell_0.5.0
## [65] pillar_1.9.0	rappdirs_0.3.3
## [67] htmltools_0.5.7	GenomeInfoDbData_1.2.11
## [69] R6_2.5.1	dbplyr_2.4.0
## [71] rprojroot_2.0.4	lattice_0.22-5
## [73] evaluate_0.23	png_0.1-8
## [75] Rsamtools_2.18.0	memoise_2.0.1
## [77] SparseArray_1.2.4	xfun_0.42
## [79] MatrixGenerics_1.14.0	pkgconfig_2.0.3