# Working with CYCLeR

### Stefan Stefanov 9/25/2020

```
    Installation of CYCLeR
```

```
    Command line tools needed

    R packages installation

    Pre-processing the data

    Mapping with STAR

    Processing the BAM info in R

    Selecting a BSJ set

    Transcript assembly

    BSJ loci extraction

    Feature identification with SGSeq
```

 Re-couting and Processing the features Transcript prediction Output and Ouantification

> CYCLeR transcript output CYCLeR quantification

CYCLER is a pipeline for reconstruction of circRNA transcripts from RNA-seq data and their subsequent quantification. The algorithm relies on comparison between control total RNA-seq samples and circRNA enriched samples to identify circRNA specific features. Then the selected circRNA features are used to infer the transcripts through a graph-based algorithm. Once the predicted transcript set is assembled, the transcript abundances are estimated through an EM algorithm with kallisto [1]. CYCLeR takes as an input BAM files and back-splice junction (BSJ) files and outputs transcript infomation in different formats and a transcript abundance file.

## Installation of CYCLeR

## Command line tools needed

```
**STAR** - https://github.com/alexdobin/STAR
**samtools** - https://sourceforge.net/projects/samtools/files/samtools/
**kallisto** - http://pachterlab.github.io/kallisto/download
#We also suggest RStudio to supplement the interactive experience
**RStudio** - https://rstudio.com/products/rstudio/download/
```

```
R packages installation
library(devtools)
install_github("stiv1n/CYCLeR")
```

# Pre-processing the data

## Mapping with STAR The STAR [2] mapping parameters are up to a personal preference. It is imperative to include the intronMotif tag. My suggestion would be 2-pass

mapping with these parameters:

```
STAR --alignSJoverhangMin 5 --outSAMstrandField intronMotif --outFilterMismatchNmax 3
 --outFilterMismatchNoverLmax 0.1
 --chimSegmentMin 15 --chimScoreMin 1 --chimJunctionOverhangMin 15
 --outFilterMultimapNmax 50 --alignIntronMax 100000 --alignIntronMin 15
 STAR --outSAMstrandField intronMotif --outFilterMismatchNmax 7 --outFilterMismatchNoverLmax 0.3
 --alignSJoverhangMin 15 --alignSJDBoverhangMin 3
  --chimSegmentMin 15 --chimScoreMin 1 --chimJunctionOverhangMin 15
 --outFilterMultimapNmax 50 --alignIntronMax 100000 --alignIntronMin 15
 #converting to sorted BAM
 samtools view -u -h Aligned.out.sam | samtools sort <name>_sorted.bam
Processing the BAM info in R
```

## We need the information for read length, fragment length and library sizes from the BAM files.

#load the BSJ files bam\_file\_prefix<-system.file("extdata", package = "CYCLeR")</pre>

## 3 SRR1191327 enriched /home/sstefan/data/wt\_vs\_RR/SRR1191327\_sorted.bam ## 4 SRR1191335 enriched /home/sstefan/data/wt\_vs\_RR/SRR1191335\_sorted.bam

filenames<-c("sample1\_75", "sample2\_75", "sample3\_75", "sample4\_75")

BSJ\_files\_ciri<-paste0(BSJ\_files\_prefix,"/",filenames)</pre> bam\_files<-paste0(bam\_file\_prefix,"/",filenames,".bam")</pre>

```
#mark the samples control and enriched or bare the consequences
 sample_table<-data.frame(filenames,c("control","control","enriched","enriched"),bam_files,stringsAsFactors = F)</pre>
 colnames(sample_table)<-c("sample_name", "treatment", "file_bam")</pre>
 si<- DataFrame(sample_table[,c("sample_name","file_bam")])</pre>
 si$file_bam <-BamFileList(si$file_bam, asMates = T)</pre>
 #this holds all the needed info of the bam files for downstream processing
 sc <- getBamInfo(si)</pre>
 sample_table$lib_size<-sc@listData$lib_size</pre>
 sample_table$read_len<-sc@listData$read_length</pre>
Use the provided sample table template.
 ## sample_name treatment
                                                                         file_bam
 ## 1 SRR1191323 control /home/sstefan/data/wt_vs_RR/SRR1191323_sorted.bam
 ## 2 SRR1191331 control /home/sstefan/data/wt_vs_RR/SRR1191331_sorted.bam
```

### ## 2 15068605 ## 3 17144585 101 ## 4 18667943

Selecting a BSJ set

## lib\_size read\_len

## 1 19255420

Selecting a BSJ set is very important, because the algorithm assumes that the provided set of BSJ is correct. I suggest BSJ identification with CIRI2 [3] and CIRCexplorer2 [4], but the choice is up to a personal preference. I have provided some useful functions for parsing the output from BSJ identification software.

```
#load the BSJ files
BSJ_files_prefix<-paste0(system.file("extdata", package = "CYCLeR"), "/ciri_")
ciri_table<-parse.files(sample_table$sample_name,BSJ_files_prefix,"CIRI")</pre>
colnames(ciri_table)<-c("circ_id", "sample1_75", "sample2_75", "sample3_75", "sample4_75")</pre>
ciri_bsjs<-process.BSJs(ciri_table, sample_table)</pre>
# i would suggest combine the output of pipelines using different mapping tools
BSJ_files_prefix_CE<-paste0(system.file("extdata", package = "CYCLeR"), "/CE_")
ce_table<-parse.files(sample_table$sample_name,BSJ_files_prefix_CE, "CE")</pre>
colnames(ce_table)<-c("circ_id", "sample1_75", "sample2_75", "sample3_75", "sample4_75")</pre>
ce_bsjs<-process.BSJs(ce_table, sample_table)</pre>
#we need to unify the results from the BSJ identification and counting
table_circ<-combine.two.BSJ.tables(ce_bsjs,ciri_bsjs)</pre>
#further downstream we need just the mean values for enriched samples
table_circ<-table_circ[,c("chr", "start", "end", "meanRR")]</pre>
colnames(table_circ)<-c("chr", "start", "end", "count")</pre>
BSJ_set<-union(ciri_bsjs$circ_id,ce_bsjs$circ_id)
BSJ_set<-BSJ_set[!grepl("caffold", BSJ_set)]</pre>
#just in case
BSJ_set<-BSJ_set[!grepl("mitochondrion", BSJ_set)]</pre>
#converting the BSJ set into a GRanges object
BSJ_gr<-make.BSJ.gr(BSJ_set)</pre>
```

head(table\_circ) ## # A tibble: 6 x 4

The parse files can work with CIRI2, CIRCexplorer2 or TSV file. Naturally a person may have different criterion for correct BSJs based on different

criteria. It is not an issue as long as the data is presented in the following template:

```
chr start
     <chr> <chr>
                  <chr> <dbl>
         10036818 10037279 1.18
 ## 2 2L
         10080375 10081946 0.923
 ## 3 2L
         10096914 10098010 0.332
 ## 4 2L 10096914 10137559 0.292
 ## 5 2L 10104792 10114895 0.268
 ## 6 2L 10120720 10122038 0.260
Transcript assembly
```

### Prior to the feature detection the files need to be trimmed to speed up the process. Afterwards the transcript features (e.g. exons, junctions) are identified with SGSeq [5]. The files are convereted with samtools [6].

**BSJ** loci extraction

#get the gene/transcript info #restoreSeqlevels(txdb)

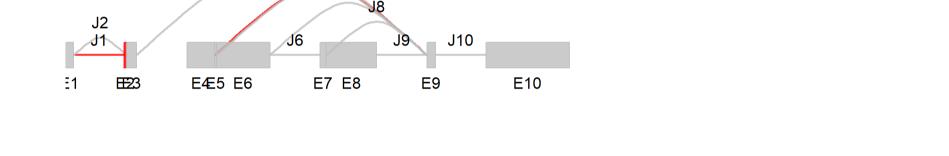
```
txdb <- TxDb.Dmelanogaster.UCSC.dm6.ensGene
 txdb <- keepSeqlevels(txdb, c("chr2L", "chr2R", "chr3R", "chr3L", "chr4", "chrX", "chrY"))
 seqlevelsStyle(txdb) <- "Ensembl"</pre>
 gene_ranges <- genes(txdb)</pre>
 txf <- convertToTxFeatures(txdb)</pre>
 #asnnotation as sg-object
 sqf <- convertToSGFeatures(txf)</pre>
 samtools_prefix<-"/home/sstefan/software/samtools-1.10/bin/"
 trimmed_bams<-filter.bam(BSJ_gr, sample_table, samtools_prefix)</pre>
 sc@listData[["file_bam"]]<-trimmed_bams</pre>
 #sc@listData[["sample_name"]]<-sample_table_ce$sample_name</pre>
Feature identification with SGSeq
 sgfc_pred <- analyzeFeatures(sc, min_junction_count=2, beta =0.1 , min_n_sample=1,cores=1,verbose=F)</pre>
```

# SGSeq feature plotting function can be used for visual representation of the control VS enriched difference

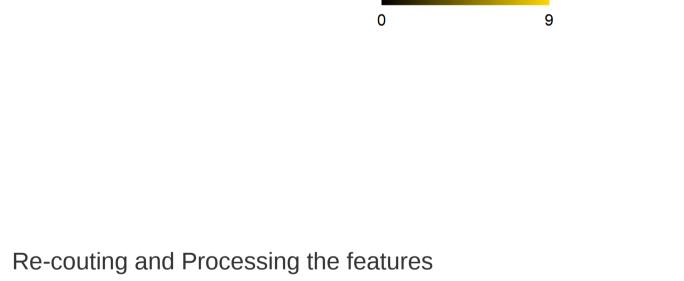
sgfc\_pred <- SGSeq::annotate(sgfc\_pred, txf)</pre>

```
quare=T)
                               J3
```

plotFeatures(sgfc\_pred, geneID = "1",assay = "counts",color\_novel = "red", include = "both",tx\_view=F,Rowv=NA, s



SRR1191323 SRR1191331 SRR1191327 SRR1191335



full\_sg<-overlap.SG.BSJ(sgfc\_pred,BSJ\_gr) #includes linear and circular features

### full\_sg<-full\_sg[rowSums(as.data.frame(full\_fc[, sample\_table\$treatment=="enriched"]))>15] full\_fc<-full\_fc[rowSums(as.data.frame(full\_fc[,sample\_table\$treatment=="enriched"]))>15,] circ\_sg<-full\_sg[full\_sg%over%BSJ\_gr] #includes features within BSJ enclosed region</pre> lin\_sg<-full\_sg[full\_sg%outside%BSJ\_gr] #includes features outside of BSJ enclosed region #annotate the BSJ with the corresponding geneIDs

# we have made new feature set so we need to recount

full\_fc<-recount.features(full\_sg, sample\_table)#fc==feature counts</pre>

simple comparison of normalized coverage values.

#removing super low coverage features

seqs<-get.seqs(full\_sg,bs\_genome)</pre>

#extract BSJ-corrected splice graphs (sg)

BSJ\_sg<-make.BSJ.sg(circ\_sg,BSJ\_gr)</pre> #full\_fc<-count\_matrix[full\_sg@featureID,]</pre> #get the correct genome for sequence info bs\_genome=Dmelanogaster #RPKM calculation for exons

full\_rpkm<-RPKM.calc(full\_fc, full\_sg, BSJ\_gr, bs\_genome=bs\_genome , sample\_table=sample\_table,

I prefer the **RSubread** [7] counting method, thus I re-count the identified exon features. I use the **SGseq** counted junctions. The features that are depleted in circRNA enriched samples need to be removed. CYCLeR provides 2 approaches for identifying depleted features: DEU strategy and

```
feature_type = "e", gc_correction = T)
 lin_rpkm<-full_rpkm[full_sg%outside%BSJ_gr,]</pre>
 #extracting circ specific counts
 circ_fc_adj<-full_rpkm[full_sg%over%BSJ_gr,]</pre>
 depleted_exons<-find.depleted.features(circ_fc_adj, sample_table, circ_sg)</pre>
 #making sure that the circ edge exons remian in the mix; they coudl be depleted in case of very low levels of the
 circle
 edge_features<-union(full_sg@featureID[start(full_sg)%in%start(BSJ_gr)],</pre>
                     full_sg@featureID[end(full_sg)%in%end(BSJ_gr)])
 depleted_exons<-setdiff(depleted_exons,edge_features)</pre>
circ_exons<-circ_sg[!circ_sg@featureID%in%depleted_exons]# the final set of circRNA exons</pre>
 circ_exons_counts<-circ_fc_adj[!circ_sg@featureID%in%depleted_exons,]</pre>
 #now for junctions
 #we need to normalize the junction read counts to the exon counts
count_matrix<-as.data.frame(counts(sgfc_pred))</pre>
 count_matrix <- apply (count_matrix, c (1, 2), function (x) {(as.integer(x))})
 sg_gr<-rowRanges(sgfc_pred)</pre>
 sg_gr_j<-sg_gr[sg_gr@type=="J"]
#circ_sg_j<-sg_gr_j[sg_gr_j%over%BSJ_gr]</pre>
 circ_sg_j<-sg_gr_j[unique(queryHits(findOverlaps(sg_gr_j,BSJ_gr,type = "within")))]</pre>
count_matrix_j<-count_matrix[circ_sg_j@featureID,]</pre>
 #get the relative sequences of around a junction
 seqs_j<-paste0(seqs[match(start(circ_sg_j),end(full_sg))],</pre>
               seqs[match(end(circ_sg_j),start(full_sg))])
 #calculate the scaled read count for junction
 junc_rpkm<-RPKM.calc(count_matrix=count_matrix_j, sg=circ_sg_j, bsj_granges = BSJ_gr, sample_table = sample_tab</pre>
 le, feature_type = "j")
deplted_j<-find.depleted.features(junc_rpkm, sample_table, circ_sg_j)</pre>
 circ_junc<-circ_sg_j[!circ_sg_j@featureID%in%deplted_j]
 circ_junc_counts<-junc_rpkm[!circ_sg_j@featureID%in%deplted_j,]</pre>
 circ_junc_counts[circ_junc_counts==0]<-1</pre>
 colnames(circ_junc_counts)<-sample_table$sample_name</pre>
The circRNA exon features are stored in SGRanges format with a corresponding matrix
 circ_exons[geneID(circ_exons)==1]
 ## SGFeatures object with 2 ranges and 0 metadata columns:
                    ranges strand type splice5p splice3p featureID
         seqnames
            <Rle> <IRanges> <Rle> <factor> <logical> <logical> <integer>
 ##
 ##
           2L 74903-75018 + E FALSE
     [1]
                                                        FALSE 70542
               2L 75078-75366 +
                                       E FALSE FALSE 70543
 ##
 ##
            geneID
                                                  txName
                                                            geneName
                                       <CharacterList> <CharacterList>
 ##
         <integer>
                1 FBtr0306540,FBtr0078101,FBtr0302164,...
 ##
     [1]
                                                             FBgn0031213
                 1 FBtr0306540,FBtr0078101,FBtr0302164,...
 ##
     [2]
                                                             FBgn0031213
     seqinfo: 1870 sequences from an unspecified genome
 circ_exons_counts[geneID(circ_exons)==1,]
```

```
Transcript prediction is processed one samples at a time. The transcript sets from different samples are then merged.
 qics_out1<-transcripts.per.sample(<sample3>)
 qics_out2<-transcripts.per.sample(<sample4>)
 qics_out_final<-merge.qics(qics_out1, qics_out2)</pre>
```

## CYCLeR provides 3 forms of output of the annotated transcript: a comprehensive flat file, a GTF-like file, and FASTA file. gtf.table<-prep.output.gtf(qics\_out\_final,circ\_exons)</pre> write.table(qics\_out\_final[,-9],file = "dm\_circles.txt", sep = "\t",row.names = F, col.names = T,quote=F)

#Kallisto comands

doi:10.1038/nbt.3519.

ONE. 2016;11:1-8.

2009;25:2078-9.

reads. Nucleic Acids Research. 2019;47.

CYCLeR transcript output

**Output and Quantification** 

qics\_out\_fa<-DNAStringSet(qics\_out\_final\$seq)</pre> names(qics\_out\_fa)<-qics\_out\_final\$circID</pre>

Transcript prediction

##

## 70542 ## 70543

#if you have a known set of circRNA in FASTA format the CYCLeR output can be combined with it fasta\_circ<-readDNAStringSet("...")</pre> final\_ref\_fa<-merge.fasta(qics\_out\_fa, fasta\_circ)</pre> writeXStringSet(qics\_out\_fa, '...')

cat linear\_transcripts.fa circles\_seq\_extended\_padded.fa > for\_kallisto.fa

kallisto quant -i kallisto\_index -o ./ sample\_1.fastq sample\_2.fastq

kallisto index -i kallisto\_index -k 31 for\_kallisto.fa

SRR1191323 SRR1191331 SRR1191327 SRR1191335

```
CYCLeR quantification
The final transcript abundance estimation is performed with kallisto. An extended and padded circRNA reference sequences are build and
combined with linear RNA sequences kallisto index is created to be used for any desired sample quantification.
 #extanding the sequence
 d.fa
 #Padding the fasta
 perl -pe 's/^[^>].*/"N"x159 . "$&"/e' circles_seq_extended.fa > circles_seq_extended_padded.fa
 #merging linear and circular sequences
```

```
2. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: Ultrafast universal RNA-seq aligner. Bioinformatics. 2013;29:15–
21.
3. Gao Y, Zhang J, Zhao F. Circular RNA identification based on multiple seed matching. Briefings in bioinformatics. 2018;19:803–10.
4. Zhang X, Dong R, Zhang Y, Zhang J, Luo Z, Zhang J, et al. Diverse alternative back-splicing and alternative splicing landscape of circular RNAs.
Genome Research. 2016;1277–87.
```

5. Goldstein LD, Cao Y, Pau G, Lawrence M, Wu TD, Seshagiri S, et al. Prediction and quantification of splice events from RNA-seq data. PLoS

7. Liao Y, Smyth GK, Shi W. The R package Rsubread is easier, faster, cheaper and better for alignment and quantification of RNA sequencing

6. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map format and SAMtools. Bioinformatics.

1. Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq quantification. Nature Biotechnology. 2016;34:525–7.