

PRAM: Pooling RNA-seq and Assembling Models

Peng Liu, Colin N. Dewey, and Sündüz Keleş

2019-10-29

Contents

1	Introduction	3
2	Installation	3
2.1	From GitHub	3
2.2	From Bioconductor	3
3	Quick start.	4
3.1	Description	4
3.2	Examples	4
4	Define intergenic genomic ranges: <code>defIgRanges()</code>	5
4.1	Description	5
4.2	Example	5
5	Prepare input RNA-seq alignments: <code>prepIgBam()</code>	6
5.1	Description	6
5.2	Example	6
6	Build transcript models: <code>buildModel()</code>	6
6.1	Description	6
6.2	Transcript prediction methods	7
6.3	Required external software	7
6.4	Example	7
7	Select transcript models: <code>selModel()</code>	8
7.1	Description	8
7.2	Example	8
8	Evaluate transcript models: <code>evalModel()</code>	9
8.1	Motivation	9

PRAM: Pooling RNA-seq and Assembling Models

8.2	Input	9
8.3	Output	9
8.4	Example	10
9	Session Info	10

1 Introduction

Pooling RNA-seq and Assembling Models (**PRAM**) is an **R** package that utilizes multiple RNA-seq datasets to predict transcript models. The workflow of PRAM contains four steps. Figure 1 shows each step with function name and associated key parameters. In addition, we provide a function named `evalModel()` to evaluate prediction accuracy by comparing transcript models with true transcripts. In the later sections of this vignette, we will describe each function in details.

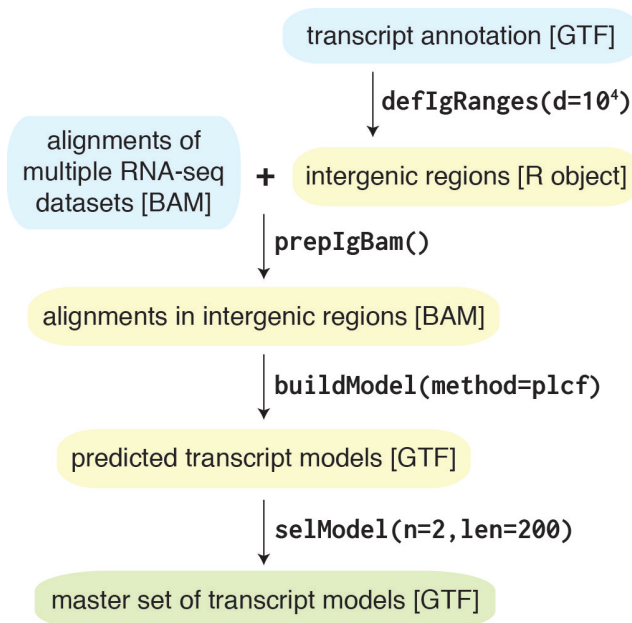


Figure 1: **PRAM** workflow

2 Installation

2.1 From GitHub

Start **R** and enter:

```
devtools::install_github('pliu55/pram')
```

2.2 From Bioconductor

Start **R** and enter:

```
if (!requireNamespace("BiocManager", quietly = TRUE))
  install.packages("BiocManager")
BiocManager::install("pram")
```

For different versions of R, please refer to the appropriate [Bioconductor release](#).

3 Quick start

3.1 Description

PRAM provides a function named `runPRAM()` to let you conveniently run through the whole workflow.

For a given gene annotation and RNA-seq alignments, you can predict transcript models in intergenic genomic regions:

```
##
## assuming the stringtie binary is in folder /usr/local/stringtie-1.3.3/
##
pram::runPRAM( in_gtf, in_bamv, out_gtf, method='plst',
               stringtie='/usr/loca/stringtie-1.3.3/stringtie')
```

- `in_gtf`: an input GTF file defining genomic coordinates of existing genes. Required to have an attribute of `gene_id` in the ninth column.
- `in_bamv`: a vector of input BAM file(s) containing RNA-seq alignments. Currently, PRAM only supports strand-specific paired-end RNA-seq with the first mate on the right-most of transcript coordinate, i.e., 'fr-firststrand' by Cufflinks definition.
- `out_gtf`: an output GTF file of predicted transcript models.
- `method`: prediction method. For the command above, we were using pooling RNA-seq datasets and building models by **stringtie**. For a list of available PRAM methods, please check Table 2 below.
- `stringtie`: location of the **stringtie** binary. PRAM's model-building step depends on external software. For more information on this topic, please see Section 6.3 below.

3.2 Examples

PRAM has included input examples files in its `extdata/demo/` folder. The table below provides a quick summary of all the example files.

Table 1: `runPRAM()`'s input example files

input argument	file name(s)
<code>in_gtf</code>	<code>in.gtf</code>
<code>in_bamv</code>	<code>SZP.bam</code> , <code>TLC.bam</code>

You can access example files by `system.file()` in R, e.g. for the argument `in_gtf`, you can access its example file by

```
system.file('extdata/demo/in.gtf', package='pram')
#> [1] "/tmp/RtmpC6sCXe/Rinst10fe12f380fe/pram/extdata/demo/in.gtf"
```

Below shows the usage of `runPRAM()` with example input files:

```
in_gtf = system.file('extdata/demo/in.gtf', package='pram')

in_bamv = c(system.file('extdata/demo/SZP.bam', package='pram'),
            system.file('extdata/demo/TLC.bam', package='pram') )

pred_out_gtf = tempfile(fileext='.gtf')

##
## assuming the stringtie binary is in folder /usr/local/stringtie-1.3.3/
##
pram::runPRAM( in_gtf, in_bamv, pred_out_gtf, method='plst',
               stringtie='/usr/loca/stringtie-1.3.3/stringtie')
```

4 Define intergenic genomic ranges: `defIgRanges()`

4.1 Description

To predict intergenic transcripts, we must first define intergenic regions by `defIgRanges()`. This function requires a GTF file containing known gene annotation supplied for its `in_gtf` argument. This GTF file should contain an attribute of **gene_id** in its ninth column. We provided an example input GTF file in PRAM package: `extdata/gtf/defIgRanges_in.gtf`.

In addition to gene annotation, `defIgRanges()` also requires user to provide chromosome sizes so that it would know the maximum genomic ranges. You can provide one of the following arguments:

- `chromgrs`: a GRanges object, or
- `genome`: a genome name, currently supported ones are: **hg19**, **hg38**, **mm9**, and **mm10**, or
- `fchromsize`: a UCSC genome browser-style size file, e.g. **hg19**

By default, `defIgRanges()` will define intergenic ranges as regions 10 kb away from any known genes. You can change it by the `radius` argument.

4.2 Example

```
pram::defIgRanges(system.file('extdata/gtf/defIgRanges_in.gtf', package='pram'),
                  genome = 'hg38')

#> GRanges object with 456 ranges and 0 metadata columns:
#>
#>      seqnames      ranges strand
#>      <Rle>        <IRanges>  <Rle>
#> [1]      chr1 30001-248956422      *
#> [2]      chr2   1-242193529      *
#> [3]      chr3   1-198295559      *
#> [4]      chr4   1-190214555      *
#> [5]      chr5   1-181538259      *
#> ...      ...      ...      ...
```

```
#> [452] chrUn_KI270753v1      1-62944      *
#> [453] chrUn_KI270754v1      1-40191      *
#> [454] chrUn_KI270755v1      1-36723      *
#> [455] chrUn_KI270756v1      1-79590      *
#> [456] chrUn_KI270757v1      1-71251      *
#> -----
#> seqinfo: 455 sequences from an unspecified genome; no seqlengths
```

5 Prepare input RNA-seq alignments: `prepIgBam()`

5.1 Description

Once intergenic regions were defined, `prepIgBam()` will extract corresponding RNA-seq alignments from input BAM files. In this way, transcript models predicted at later stage will solely from intergenic regions. Also, with fewer RNA-seq alignments, model prediction will run faster.

Three input arguments are required by `prepIgBam()`:

- `finbam`: an input RNA-seq BAM file sorted by genomic coordinate. Currently, we only support strand-specific paired-end RNA-seq data with the first mate on the right-most of transcript coordinate, i.e. 'fr-firststrand' by Cufflinks's definition.
- `iggrs`: a `GRanges` object to define intergenic regions.
- `foutbam`: an output BAM file.

5.2 Example

```
finbam = system.file('extdata/bam/CMPRep2.sortedByCoord.raw.bam',
                     package='pram')

iggrs = GenomicRanges::GRanges('chr10:77236000-77247000:+')

foutbam = tempfile(fileext='.bam')

pram::prepIgBam(finbam, iggrs, foutbam)
#> Extracted RNA-seq alignments are saved in /tmp/Rtmp0qrFx7/file140b578f51a2.bam
```

6 Build transcript models: `buildModel()`

6.1 Description

`buildModel()` predict transcript models from RNA-seq BAM file(s). This function requires two arguments:

- `in_bamv`: a vector of input BAM file(s)

- `out_gtf`: an output GTF file containing predicted transcript models

6.2 Transcript prediction methods

`buildModel()` has implemented seven transcript prediction methods. You can specify it by the `method` argument with one of the keywords: **plcf**, **plst**, **cfmg**, **cftc**, **stmg**, **cf**, and **st**. The first five denote meta-assembly methods that utilize multiple RNA-seq datasets to predict a single set of transcript models. The last two represent methods that predict transcript models from a single RNA-seq dataset.

The table below compares prediction steps for these seven methods. By default, `buildModel()` uses **plcf** to predict transcript models.

Table 2: Prediction steps of the seven `buildModel()` methods

method	meta-assembly	preparing RNA-seq input	building transcripts	assembling transcripts
plcf	yes	pooling alignments	Cufflinks	no
plst	yes	pooling alignments	StringTie	no
cfmg	yes	no	Cufflinks	Cuffmerge
cftc	yes	no	Cufflinks	TACO
stmg	yes	no	StringTie	StringTie-merge
cf	no	no	Cufflinks	no
st	no	no	StringTie	no

6.3 Required external software

Depending on your specified prediction method, `buildModel()` requires external software: Cufflinks, StringTie and/or TACO, to build and/or assemble transcript models. You can either specify the software location using the `cufflinks`, `stringtie`, and `taco` arguments in `buildModel()`, or simply leave these three arguments undefined and let PRAM download them for you automatically. The table below summarized software versions `buildModel()` would download when required software was not specified. Please note that, for **macOS**, pre-compiled Cufflinks binary versions 2.2.1 and 2.2.0 appear to have an issue on processing BAM files, therefore we recommend to use version 2.1.1 instead.

Table 3: `buildModel()`-required software and recommended version

software	Linux binary	macOS binary	required by
Cufflinks, Cuffmerge	v2.2.1	v2.1.1	plcf , cfmg , cftc , and cf plst , stmg , and st cftc
StringTie, StringTie-merge	v1.3.3b	v1.3.3b	
TACO	v0.7.0	v0.7.0	

6.4 Example

```
fbams = c( system.file('extdata/bam/CMPRep1.sortedByCoord.clean.bam',
                      package='pram'),
           system.file('extdata/bam/CMPRep2.sortedByCoord.clean.bam',
                      package='pram') )

foutgtf = tempfile(fileext='.gtf')

##
## assuming the stringtie binary is in folder /usr/local/stringtie-1.3.3/
##
pram::buildModel(fbams, foutgtf, method='plst',
                 stringtie='/usr/local/stringtie-1.3.3/stringtie')
```

7 Select transcript models: `selModel()`

7.1 Description

Once transcript models were built, you may want to select a subset of them by their genomic features. `selModel()` was developed for this purpose. It allows you to select transcript models by their total number of exons and total length of exons and introns.

`selModel()` requires two arguments:

- `fin_gtf`: input GTF file containing to-be-selected transcript models. This file is required to have **transcript_id** attribute in the ninth column.
- `fout_gtf`: output GTF file containing selected transcript models.

By default: `selModel()` will select transcript models with ≥ 2 exons and ≥ 200 bp total length of exons and introns. You can change the default using the `min_n_exon` and `min_tr_len` arguments.

7.2 Example

```
fin_gtf = system.file('extdata/gtf/selModel_in.gtf', package='pram')

fout_gtf = tempfile(fileext='.gtf')

pram::selModel(fin_gtf, fout_gtf)
#> Selected transcript models are saved in /tmp/Rtmp0qrFx7/file140b2b48d4.gtf
```

8 Evaluate transcript models: `evalModel()`

8.1 Motivation

After PRAM has predicted a number of transcript models, you may wonder how accurate these models are. To answer this question, you can compare PRAM models with real transcripts (i.e., positive controls) that you know should be predicted. PRAM's `evalModel()` function will help you to make such comparison. It will calculate precision and recall rates on three features of a transcript: exon nucleotides, individual splice junctions, and transcript structure (i.e., whether all splice junctions within a transcript were constructed in a model).

8.2 Input

`evalModel()` requires two arguments:

- `model_exons`: genomic coordinates of transcript model exons.
- `target_exons`: genomic coordinates of real transcript exons.

The two arguments can be in multiple formats:

- both are `GRanges` objects
- both are `character` objects denoting names of GTF files
- both are `data.table` objects containing the following five columns for each exon:
 - **chrom**: chromosome name
 - **start**: starting coordinate
 - **end**: ending coordinate
 - **strand**: strand information, either '+' or '-'
 - **trid**: transcript ID
- `model_exons` is the name of a GTF file and `target_exons` is a `data.table` object.

8.3 Output

The output of `evalModel()` is a `data.table` object, where columns are evaluation results and each row is three transcript features.

Table 4: `evalModel()` output columns

column name	representation
feat	transcript feature
ntp	number of true positives (TP)
nfn	number of false negatives (FN)
nfp	number of false positives (FP)
precision	precision rate: $\frac{TP}{(TP+FP)}$
recall	recall rate: $\frac{TP}{(TP+FN)}$

Table 5: `evalModel()` output rows

feature name	representation
exon_nuc	exon nucleotide
indi_jnc	individual splice junction
tr_jnc	transcript structure

8.4 Example

```
fmdl = system.file('extdata/benchmark/plcf.tsv', package='pram')
ftgt = system.file('extdata/benchmark/tgt.tsv', package='pram')

mdltdt = data.table::fread(fmdl, header=TRUE, sep="\t")
tgttdt = data.table::fread(ftgt, header=TRUE, sep="\t")

pram::evalModel(mdltdt, tgttdt)
#>      feat      ntp      nfn      nfp precision      recall
#> 1: exon_nuc 1723581 109337 8424 0.9951363 0.9403481
#> 2: indi_jnc   2889    162   192 0.9376826 0.9469027
#> 3:  tr_jnc   1138    118   252 0.8187050 0.9060510
```

9 Session Info

Below is the output of `sessionInfo()` on the system on which this document was compiled.

```
#> R Under development (unstable) (2019-10-24 r77329)
#> Platform: x86_64-pc-linux-gnu (64-bit)
#> Running under: Ubuntu 18.04.3 LTS
#>
#> Matrix products: default
#> BLAS: /home/biocbuild/bbs-3.11-bioc/R/lib/libRblas.so
#> LAPACK: /home/biocbuild/bbs-3.11-bioc/R/lib/libRlapack.so
#>
#> locale:
#>  [1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
#>  [3] LC_TIME=en_US.UTF-8      LC_COLLATE=C
#>  [5] LC_MONETARY=en_US.UTF-8  LC_MESSAGES=en_US.UTF-8
#>  [7] LC_PAPER=en_US.UTF-8     LC_NAME=C
#>  [9] LC_ADDRESS=C             LC_TELEPHONE=C
#> [11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
#>
#> attached base packages:
#> [1] stats      graphics  grDevices  utils      datasets  methods   base
#>
#> other attached packages:
#> [1] BiocStyle_2.15.0
#>
```

PRAM: Pooling RNA-seq and Assembling Models

```
#> loaded via a namespace (and not attached):
#> [1] Rcpp_1.0.2                compiler_4.0.0
#> [3] BiocManager_1.30.9        GenomeInfoDb_1.23.0
#> [5] XVector_0.27.0            bitops_1.0-6
#> [7] tools_4.0.0               zlibbioc_1.33.0
#> [9] digest_0.6.22             evaluate_0.14
#> [11] lattice_0.20-38           rlang_0.4.1
#> [13] Matrix_1.2-17             DelayedArray_0.13.0
#> [15] yaml_2.2.0                parallel_4.0.0
#> [17] xfun_0.10                 GenomeInfoDbData_1.2.2
#> [19] rtracklayer_1.47.0        stringr_1.4.0
#> [21] knitr_1.25                Biostrings_2.55.0
#> [23] S4Vectors_0.25.0          IRanges_2.21.0
#> [25] stats4_4.0.0              grid_4.0.0
#> [27] data.table_1.12.6         Biobase_2.47.0
#> [29] XML_3.98-1.20             BiocParallel_1.21.0
#> [31] rmarkdown_1.16            bookdown_0.14
#> [33] magrittr_1.5              Rsamtools_2.3.0
#> [35] htmltools_0.4.0           matrixStats_0.55.0
#> [37] GenomicAlignments_1.23.0  BiocGenerics_0.33.0
#> [39] GenomicRanges_1.39.0     SummarizedExperiment_1.17.0
#> [41] stringi_1.4.3             RCurl_1.95-4.12
#> [43] pram_1.3.0
```