

How to use breakpointR

*David Porubsky**

*david.porubsky@gmail.com

October 30, 2018

Contents

1	Introduction	2
2	Quickstart	2
2.1	Running breakpointR	3
3	Recommended settings	4
3.1	Reading BAM files	4
3.2	Removing certain regions	4
3.3	Binning strategy	4
3.4	Breakpoint peak detection	4
3.5	Background reads	5
3.6	Calling breakpoint hotspots	5
3.7	Loading results and plotting single cells	6
4	Session Info	7

1 Introduction

BreakpointR is a novel algorithm designed to accurately tracks template strand changes in Strand-seq data using a bi-directional read-based binning. Read-based binning strategy scales bin size dynamically to accommodate defined number of reads what accounts for mappability bias in sparsely covered single-cell Strand-seq data. In such dynamically scaled bins, read directionality is tracked in order to search for points where template strand state changes. BreakpointR takes as an input reads aligned to the reference genome stored in BAM files. BreakpointR outputs locations where directionality of sequenced teplate strands changes.

2 Quickstart

The main function of this package is called `breakpointR()` and performs all the necessary steps to get from aligned reads in BAMs to interpretable output. For an unexperienced user we advise to run `breakpointR` with default parameters and later based on the obtained results start to tweak certain parameters. For more detailed guidance on parameter tweaking see section 3.

```
library(breakpointR)
## Run breakpointR with default paprameters
breakpointR(inputfolder='folder-with-BAM', outputfolder='output-directory')
```

Although in most cases the one of the above commands will produce reasonably good results, it might be worthwile to adjust the default parameters to improve performance and the quality of the results. You can get a description of all available parameters by typing

```
?breakpointR
```

After the function has finished, you will find the folder **output-directory** containing all produced files and plots. This folder contains the following **files** and **folders**:

- `breakpointR.config`: This file contains all the parameters that are necessary to reproduce your analysis. You can specify this file as

```
breakpointR(..., configfile='breakpointR.config')
```

to run another analysis with the same parameter settings.

- **breakpoints** UCSC browser formatted bedgraphs compiling all breakpoints across all single-cell libraries. This folder also contains list of all localized breakpoints in all single-cell libraries. Lastly, locations of breakpoint hotspots are reported here if

```
callHotSpots=TRUE
```

- **browserfiles** UCSC browser formatted files with exported reads, deltaWs and break-Points for every single-cell library.
- **data** Contains RData files storing results of BreakpointR analysis for each single-cell library.

How to use breakpointR

- **plots:** Genome-wide plots for selected chromosome, genome-wide heatmap of strand states as well as chromosome specific read distribution together with localized breakpoints. All plots are created by default.

2.1 Running breakpointR

The function `breakpointR()` takes an input BAM files stored in the inputfolder and produces an output folder with results, plots and browserfiles. The following code is an example of how to run `breakpointR` for single-end reads with window size defined by size. Results will be stored in **outputfolder/data** as RData objects. Such data can be later loaded for further processing and customized plotting.

```
library(breakpointR)

## Get some example files
datafolder <- system.file("extdata", "example_bams", package="breakpointRdata")
outputfolder <- tempdir()
## Run breakpointR
breakpointR(inputfolder = datafolder, outputfolder = outputfolder,
            chromosomes = 'chr22', pairedEndReads = FALSE,
            reuse.existing.files = FALSE, window size = 1000000,
            binMethod = 'size', pair2frgm = FALSE, min.mapq = 10,
            filtAlt = TRUE)
```

3 Recommended settings

3.1 Reading BAM files

Currently *breakpointR* can take as an input only aligned reads stored in BAM files. All BAM files are expected to be present in a folder specified in `breakpointR(..., inputfolder)`. We advise to remove reads with low mapping quality and reads with alternative alignments. Duplicated reads are removed by default (to keep them).

```
breakpointR(..., min.mapq = 10, filtAlt = TRUE)
```

3.2 Removing certain regions

breakpointR allows a user to exclude certain genomic regions from the analysis. This comes handy when one wants to remove reads that falls into low complexity regions such as segmental duplications or centromeres. To mask certain genomic regions user has to define option `breakpointR(..., maskRegions)` to a bed formatted text file. All reads falling into these regions will be discarded prior to breakpoint detection. User defined regions to mask can be downloaded from the UCSC Table Browser.

3.3 Binning strategy

breakpointR uses read based binning strategy and offers two approaches to set the bin size: (1) user defined number of reads in each bin or (2) number of reads in every bin is selected based on desired bin length.

```
library(breakpointR)
## Binning strategy based on desired bin length
breakpointR(inputfolder='folder-with-BAM', outputfolder='output-directory',
            windowSize=1e6, binMethod='size')
## Binning strategy based user-defined number of reads in each bin
breakpointR(inputfolder='folder-with-BAM', outputfolder='output-directory',
            windowSize=100, binMethod='reads')
```

Based on the size of the user defined bin changes sensitivity and specificity of breakpoint detection. We recommend to select rather large bin size ($\geq 1\text{Mb}$) in order to reliably detect low frequency SCE events. In order to detect smaller events like inversions smaller bin size is recommended. Keep in mind that such settings also leads to higher level of false positive breakpoints. In this case one might need to tweak some breakpoint detection parameters (see subsection 3.4).

3.4 Breakpoint peak detection

Breakpoint detection is based on finding significant peaks in ΔW values. Level of significance is measured in z-score (or number SDs) from the set threshold `breakpointR(..., peakTh)`. By default the threshold is set to the 1/3 of the highest ΔW value. For the data

How to use breakpointR

with noisy coverage we recommend to set this threshold little bit higher, for example 1/2 of the highest deltaW value. In case of noisy data we also recommend to tweak trim option `breakpointR(..., trim)` which used to calculate SD after trimming extreme deltaW values.

```
## Example deltaW values
exampleFolder <- system.file("extdata", "example_results",
                             package="breakpointRdata")
exampleFile <- list.files(exampleFolder, full.names=TRUE)[1]
breakpoint.object <- loadFromFiles(exampleFile)
head(breakpoint.object[[1]]$deltas)

## GRanges object with 6 ranges and 1 metadata column:
##      seqnames      ranges strand |      deltaW
##      <Rle>      <IRanges> <Rle> | <numeric>
## [1]   chr1    7560-7594     - |         17
## [2]   chr1    8569-8612     + |          0
## [3]   chr1 15116-15143     - |         57
## [4]   chr1 17235-17240     - |        130
## [5]   chr1 19615-19720     - |         41
## [6]   chr1 19849-19911     - |         38
## -----
##      seqinfo: 23 sequences from an unspecified genome
```

3.5 Background reads

Background reads are a common feature of Strand-seq libraries. Strand-seq is based on removal of newly synthesized strand during DNA replication, however this process is not perfect. Therefore, we usually expect low abundance reads aligned in opposite direction even for purely WW or CC chromosomes. Another reason to see such artefacts is imperfect read mapping especially in repetitive regions. To remove reads falling into the repetitive regions see subsection 3.2.

3.6 Calling breakpoint hotspots

In order to find locations where breakpoints occur around the same genomic position in multiple Strand-seq libraries there is `hotspotter()`. Function can be invoked by setting corresponding parameter to 'TRUE'. It make sense to set this parameter only if there is available a reasonable number (≥ 50) of Strand-seq libraries.

```
## To run breakpoint hotspot analysis using the main breakpointR function
breakpointR(..., callHotSpots=TRUE)
```

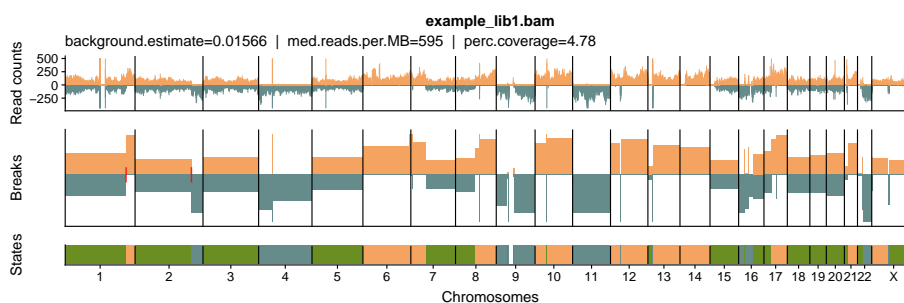
```
## To run breakpoint hotspot analysis using exported data
exampleFolder <- system.file("extdata", "example_results",
                             package="breakpointRdata")
exampleFiles <- list.files(exampleFolder, full.names=TRUE)
breakpoint.objects <- loadFromFiles(exampleFiles)
## Extract breakpoint coordinates
breaks <- lapply(breakpoint.objects, '[', 'breaks')
```

How to use breakpointR

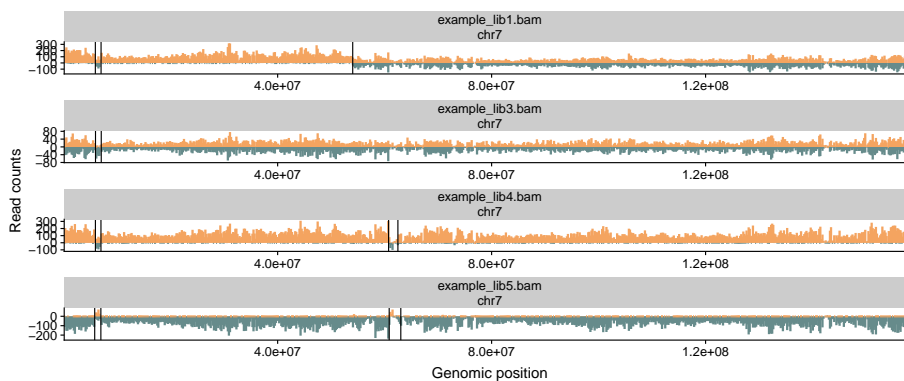
```
## Get hotspot coordinates  
hotspots <- hotspotter(breaks, bw=1e6)
```

3.7 Loading results and plotting single cells

```
## Plotting a single library  
exampleFolder <- system.file("extdata", "example_results",  
                             package="breakpointRdata")  
exampleFile <- list.files(exampleFolder, full.names=TRUE)[1]  
plotBreakpoints(exampleFile)  
  
## [[1]]
```



```
## Plotting a single library  
exampleFolder <- system.file("extdata", "example_results",  
                             package="breakpointRdata")  
exampleFiles <- list.files(exampleFolder, full.names=TRUE)[1:4]  
plotBreakpointsPerChr(exampleFiles, chromosomes = 'chr7')  
  
## $chr7
```



4 Session Info

```
toLatex(sessionInfo())
```

- R version 3.5.1 Patched (2018-07-12 r74967), x86_64-pc-linux-gnu
- Locale: LC_CTYPE=en_US.UTF-8, LC_NUMERIC=C, LC_TIME=en_US.UTF-8, LC_COLLATE=C, LC_MONETARY=en_US.UTF-8, LC_MESSAGES=en_US.UTF-8, LC_PAPER=en_US.UTF-8, LC_NAME=C, LC_ADDRESS=C, LC_TELEPHONE=C, LC_MEASUREMENT=en_US.UTF-8, LC_IDENTIFICATION=C
- Running under: Ubuntu 16.04.5 LTS
- Matrix products: default
- BLAS: /home/biocbuild/bbs-3.8-bioc/R/lib/libRblas.so
- LAPACK: /home/biocbuild/bbs-3.8-bioc/R/lib/libRlapack.so
- Base packages: base, datasets, grDevices, graphics, methods, parallel, stats, stats4, utils
- Other packages: BiocGenerics 0.28.0, GenomelInfoDb 1.18.0, GenomicRanges 1.34.0, IRanges 2.16.0, S4Vectors 0.20.0, breakpointR 1.0.0, breakpointRdata 0.99.2, cowplot 0.9.3, ggplot2 3.1.0, knitr 1.20
- Loaded via a namespace (and not attached): Biobase 2.42.0, BiocManager 1.30.3, BiocParallel 1.16.0, BiocStyle 2.10.0, Biostrings 2.50.0, DelayedArray 0.8.0, GenomelInfoDbData 1.2.0, GenomicAlignments 1.18.0, Matrix 1.2-14, R6 2.3.0, RCurl 1.95-4.11, Rcpp 0.12.19, Rsamtools 1.34.0, SummarizedExperiment 1.12.0, XVector 0.22.0, assertthat 0.2.0, backports 1.1.2, bindr 0.1.1, bindrcpp 0.2.2, bitops 1.0-6, codetools 0.2-15, colorspace 1.3-2, compiler 3.5.1, crayon 1.3.4, digest 0.6.18, doParallel 1.0.14, dplyr 0.7.7, evaluate 0.12, foreach 1.4.4, glue 1.3.0, grid 3.5.1, gtable 0.2.0, gtools 3.8.1, highr 0.7, htmltools 0.3.6, iterators 1.0.10, labeling 0.3, lattice 0.20-35, lazyeval 0.2.1, magrittr 1.5, matrixStats 0.54.0, munsell 0.5.0, pillar 1.3.0, pkgconfig 2.0.2, plyr 1.8.4, purrr 0.2.5, rlang 0.3.0.1, rmarkdown 1.10, rprojroot 1.3-2, scales 1.0.0, stringi 1.2.4, stringr 1.3.1, tibble 1.4.2, tidyselect 0.2.5, tools 3.5.1, withr 2.1.2, yaml 2.2.0, zlibbioc 1.28.0