

How to use breakpointR

David Porubsky*

*david.porubsky@gmail.com

October 24, 2023

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	5
4	Session Info	7

1 Introduction

BreakpointR is a novel algorithm designed to accurately track template strand changes in Strand-seq data using a bi-directional read-based binning. Read-based binning strategy scales each bin size dynamically to accommodate a defined number of reads, which 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 and stored in a single BAM file per single cell. BreakpointR outputs locations where directionality of sequenced template strands changes. Template strand changes are defined by changes in proportion of reads mapped to positive ('Crick') and negative ('Watson') strand of the reference genome. In a diploid organism such as human we distinguish three possible scenarios of template strand inheritance. If both parental homologues were inherited as Watson template (we assign a WW state), if only Crick templates were inherited (we assign CC state), or one Watson and one Crick template was inherited by each parent (we assign WC state).

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 predicted breakpoints (changes) in strand directionality. 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 how to tweak certain parameters see section 3.

```
library(breakpointR)
## Run breakpointR with a default parameters
breakpointR(inputfolder='folder-with-BAMs', outputfolder='output-folder')
```

Although in most cases the one of the above commands will produce reasonably good results, it might be worthwhile to adjust the default parameters in order 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 parameters that are necessary to reproduce your analysis. You can specify this file as shown below in order to run another analysis with the same parameter settings.

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

- **breakpoints** UCSC browser formatted bedgraphs compile 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 breakpoints for every single-cell library.

How to use breakpointR

- **data** Contains RData files that store complete results of breakpointR analysis for each single-cell library.
- **plots**: Genome-wide plots for selected chromosomes, genome-wide heatmap of strand states as well as chromosome specific read distribution together with localized breakpoints. All aforementioned plots are created by default.

2.1 Running breakpointR

The function `breakpointR()` takes as an input BAM files stored in the inputfolder and produces an outputfolder with results, plots and browserfiles. The following code is an example of how to run `breakpointR` on single-end reads with a 'windowsize' defined by size of 1Mb (see subsection 3.3). 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, windowsize = 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 as `breakpointR(..., inputfolder)`. We advise to remove reads with low mapping quality and reads with alternative alignments. Duplicated reads are removed by default.

```
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. Such low complexity regions might cause false positive breakpoints due to the spurious mappings of short reads. To mask certain genomic regions user has to define path to a bed formatted text file as `breakpointR(..., maskRegions)`. 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 genomic size of each bin.

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

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

3.4 Breakpoint peak detection

Breakpoint detection is based on finding significant peaks in deltaW values. Level of peak significance is measured in the number of standard deviations (SD) from the set threshold (z-score) `breakpointR(..., peakTh)`. By default the threshold is set to the 1/3 of the highest deltaW value. For the data with noisy and uneven coverage we recommend to set higher threshold, for example 1/2 of the highest deltaW value. In addition, we also recommend to tweak 'trim' option `breakpointR(..., trim)` in order to set a fraction of extreme deltaW values to be excluded from SD calculation.

How to use breakpointR

```
## 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 a opposite direction even for purely WW or CC chromosomes. An another reason to see such artifacts is imperfect read mapping especially in repetitive and complex genomic regions. To remove reads falling into such 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)
```

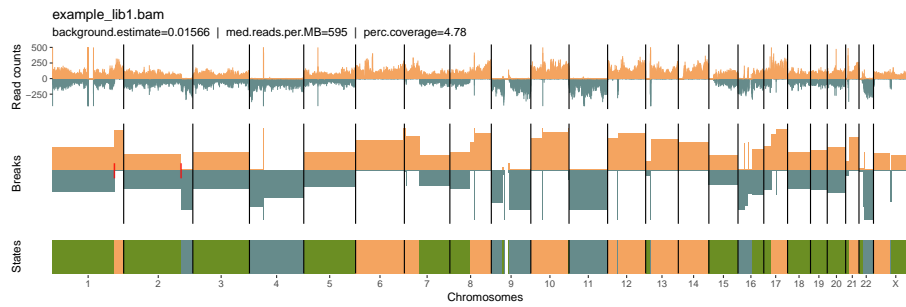
```
## To run breakpoint hotspot analysis using exported data
exampleFolder <- system.file("extdata", "example_results",
                             package="breakpointRdata")
exampleFiles <- list.files(exampleFolder, full.names=TRUE)
breakpoint.objects <- loadFromFiles(exampleFiles)
## Extract breakpoint coordinates
breaks <- lapply(breakpoint.objects, '[', 'breaks')
## Get hotspot coordinates
hotspots <- hotspotter(breaks, bw=1e6)
```

3.7 Loading results and plotting single cells

How to use breakpointR

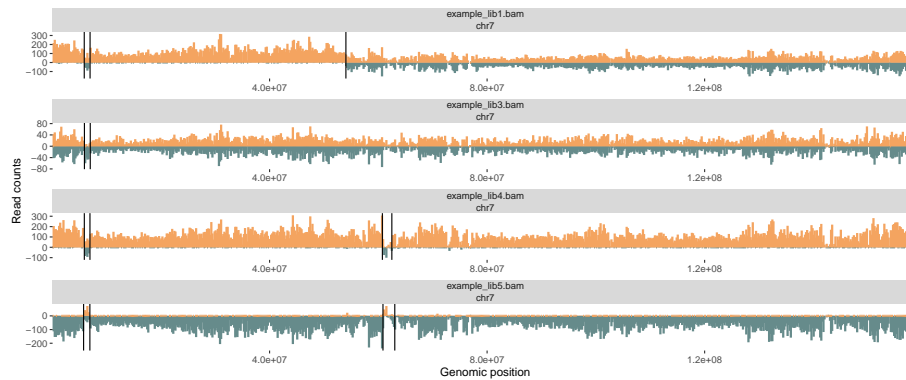
```
## Plotting a single library
exampleFolder <- system.file("extdata", "example_results",
                             package="breakpointRdata")
exampleFile <- list.files(exampleFolder, full.names=TRUE)[1]
plotBreakpoints(exampleFile)

## [[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 4.3.1 (2023-06-16), x86_64-pc-linux-gnu
- Locale: LC_CTYPE=en_US.UTF-8, LC_NUMERIC=C, LC_TIME=en_GB, 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
- Time zone: America/New_York
- TZcode source: system (glibc)
- Running under: Ubuntu 22.04.3 LTS
- Matrix products: default
- BLAS: /home/biocbuild/bbs-3.18-bioc/R/lib/libRblas.so
- LAPACK: /usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.10.0
- Base packages: base, datasets, grDevices, graphics, methods, stats, stats4, utils
- Other packages: BiocGenerics 0.48.0, GenomeInfoDb 1.38.0, GenomicRanges 1.54.0, IRanges 2.36.0, S4Vectors 0.40.0, XVector 0.42.0, breakpointR 1.20.0, breakpointRdata 1.19.0, cowplot 1.1.1, knitr 1.44
- Loaded via a namespace (and not attached): Biobase 2.62.0, BiocManager 1.30.22, BiocParallel 1.36.0, BiocStyle 2.30.0, Biostrings 2.70.0, DelayedArray 0.28.0, GenomeInfoDbData 1.2.11, GenomicAlignments 1.38.0, Matrix 1.6-1.1, MatrixGenerics 1.14.0, R6 2.5.1, RCurl 1.98-1.12, Rsamtools 2.18.0, S4Arrays 1.2.0, SparseArray 1.2.0, SummarizedExperiment 1.32.0, abind 1.4-5, bitops 1.0-7, cli 3.6.1, codetools 0.2-19, colorspace 2.1-0, compiler 4.3.1, crayon 1.5.2, digest 0.6.33, doParallel 1.0.17, dplyr 1.1.3, evaluate 0.22, fansi 1.0.5, farver 2.1.1, fastmap 1.1.1, foreach 1.5.2, generics 0.1.3, ggplot2 3.4.4, glue 1.6.2, grid 4.3.1, gtable 0.3.4, gtools 3.9.4, highr 0.10, htmltools 0.5.6.1, iterators 1.0.14, labeling 0.4.3, lattice 0.22-5, lifecycle 1.0.3, magrittr 2.0.3, matrixStats 1.0.0, munsell 0.5.0, parallel 4.3.1, pillar 1.9.0, pkgconfig 2.0.3, rlang 1.1.1, rmarkdown 2.25, scales 1.2.1, tibble 3.2.1, tidyselect 1.2.0, tools 4.3.1, utf8 1.2.4, vctrs 0.6.4, withr 2.5.1, xfun 0.40, yaml 2.3.7, zlibbioc 1.48.0