# Package 'DuplexDiscovereR'

October 24, 2025

Title Analysis of the data from RNA duplex probing experiments

**Description** DuplexDiscovereR is a package designed for analyzing data from RNA cross-linking and proximity ligation protocols such as SPLASH, PARIS, LIGR-seq, and others.

DuplexDiscovereR accepts input in the form of chimerically or split-aligned reads. It includes procedures for alignment classification, filtering, and efficient clustering of individual chimeric reads into duplex groups (DGs). Once DGs are identified, the package predicts RNA duplex formation and their hybridization energies. Additional metrics, such as p-values for random ligation hypothesis or mean DG alignment scores, can be calculated to rank final set of RNA duplexes.

Data from multiple experiments or replicates can be processed separately and further compared to check the reproducibility of the experimental method.

License GPL-3

URL https://github.com/Egors01/DuplexDiscovereR/

BugReports https://github.com/Egors01/DuplexDiscovereR/issues/

**Encoding UTF-8** 

NeedsCompilation no

Version 1.3.4

**Roxygen** list(markdown = TRUE)

RoxygenNote 7.3.2

**BiocType** Software

**biocViews** Sequencing, Transcriptomics, StructuralPrediction, Clustering, SplicedAlignment

Imports Gviz, Biostrings, rtracklayer, GenomicAlignments, GenomicRanges, ggsci, igraph, rlang, scales, stringr, dplyr, tibble, tidyr, purrr, methods, grDevices, stats, utils, vctrs

**Depends** R (>= 4.5), InteractionSet

LazyData false

**Suggests** knitr, UpSetR, BiocStyle, rmarkdown, testthat (>= 3.0.0)

2 Contents

VignetteBuilder knitr
VignetteEngine knitr
Config/testthat/edition 3
$\textbf{git\_url} \   \texttt{https://git.bioconductor.org/packages/DuplexDiscovereR}$
git_branch devel
git_last_commit 7e5c418
git_last_commit_date 2025-10-06
Repository Bioconductor 3.23
Date/Publication 2025-10-24
Author Egor Semenchenko [aut, cre, cph] (ORCID:
Maintainer Egor Semenchenko < yegor.smb@gmail.com>

# **Contents**

.addDGidsForTmpDGs
.addGeneCounts
.annotateCisTrans
.compute_clusters_comp
.DGIdToDuplexId
.getStartEndOvl
.gv_plotboxes
.gv_updatepars
$annotate GI\ldots\ldots\ldots \ \ 8$
availableDisplayPars
$calculate Ligation Pvalues \ \ldots \ 10$
$classify Two Arm Chimeras \\  \   \dots \\  \   11$
clusterDuplexGroups
collapseIdenticalReads
collapseSimilarChimeras
$collapse\_duplex\_groups  \dots  16$
col_check_rename
$compare Multiple Interactions \ . \ . \ . \ . \ . \ . \ . \ . \ . \ $
$compute GISelf Overlaps \\  \ldots \\                               $
convert_gi_to_ranges
$dd\_get\_chimeric\_reads \ \dots \ $
dd_get_chimeric_reads_stats
$dd\_get\_duplex\_groups  .  .  .  .  .  .  .  .  .  $
$dd\_get\_reads\_classes \ \dots \ \dots \ 24$
dd_get_run_stats
drawGD,DuplexTrack-method

	DuplexDiscovereR													26
	DuplexDiscovererResults	-class												27
	DuplexTrack													28
	getChimericJunctionType													31
	getRNAHybrids													32
	getSpliceJunctionChimer													33
	get_arm_a													34
	get_arm_b													34
	get_char_count_cigar													35
	get_chimeric_junctions_c	onestrand .												35
	get_colnames_and_types	_for_input												36
	makeDfFromGi													36
	makeGiFromDf													37
	preproc_chim_junction_c	out_pe												38
	preproc_chim_junction_c	out_se												38
	preproc_generic													39
	preproc_generic_gi													40
	refresh_gi													40
	RNADuplexesGeneCoun	ts												41
	RNADuplexesRawBed .													41
	RNADuplexesRawChimS	STAR												42
	RNADuplexSampleClust	Reads												42
	RNADuplexSampleDGs													43
	RNADuplexSampleGI .													44
	runDuplexDiscoPreproc													45
	runDuplexDiscoverer													46
	SampleGeneAnnoGR													49
	SampleSmallGI													50
	SampleSpliceJncGR													50
	show,DuplexDiscovererR													51
	show,DuplexTrack-metho	od												51
	subset_gi													52
	trimAroundJunction													53
	writeGiToSAMfile													54
Index														<b>56</b>
.addD	OGidsForTmpDGs Hel	per function	to ad	d ids	to the	e dup	olex g	roup	s mis	ssed	duri	ng .	glob	al
		tering				•	Ü	•					-	

# Description

Check if there are a temporary duplex records with duplex\_id, which consist of more than one read n\_reads > 1, but does not have assigned any dg\_id as the duplex group (DG) index. Creates new dg\_id if n\_reads > 1

4 .addGeneCounts

#### Usage

```
.addDGidsForTmpDGs(gi_input)
```

# **Arguments**

gi\_input GInteractions with the dg\_id, duplex\_id and n\_reads column

#### **Details**

Meant to be used in the situations when previous collapsing steps merged two or more reads to the temporary DG with duplex\_id, but global clustering has not identified any overlap between this temporary group and other duplexes, resulting in undefined dg\_id. This function looks up for these cases and creates new dg\_id for temporary DGs, marking them as the final DGs. New dg\_id values are unique and allocated sequentially after the maximum value of dg\_id

#### Value

**GInteractions** object with new dg\_id for rows with n\_reads > 1

.addGeneCounts

Helper function to add count data to metadata of GInteractions

#### **Description**

Merges the count dataframe and interactions metadata by id\_col If key is not found, in metadata throws error

#### Usage

```
.addGeneCounts(gi, df_counts, id_col = "gene_id")
```

# Arguments

gi GInteractions

df\_counts dataframe with read counts

id\_col key to use in merge

#### Value

GInteractions with added counts

.annotateCisTrans 5

.annotateCisTrans

Annotate RNA-RNA interactions as cis- and trans-

# **Description**

Annotated each entry gi object as cis, if the .A and .B arms correspond to the same feature (i.e transcript\_id or gene\_id) If the values are are equal, then annotated with value: cis = 1, If not equal or NA: cis = 0

#### Usage

```
.annotateCisTrans(gi, id_col_base = "gene_id")
```

# **Arguments**

gi GInteractions object containing two metadata columns as feature annotation

id\_col\_base base name of the feature id columns to use. Function will look for <id\_col\_base>.A

and <id\_col\_base>.B columns and compare them

#### Value

gi GInteractions object containing cis field with 0/1 values

```
.compute_clusters_comp
```

Helper function to call clustering on each component

# Description

Helper function to call clustering on each component

# Usage

```
.compute_clusters_comp(graph, index)
```

# **Details**

Call clustering on each independent graph component. Used when decompose==TRUE

### Value

sub-graph with clusters labelled with sample-wise unique ids (cluster\_group)

6 .getStartEndOvl

 $. {\tt DGIdToDuplexId}$ 

Accessor for mapping between temporary and final cluster ids

# **Description**

Accessor for mapping between temporary and final cluster ids

# Usage

```
.DGIdToDuplexId(gi)
```

# **Arguments**

gi

#### Value

tibble

.getStartEndOvl

Helper func for calculating one-side overlap of SJ and junctions of gi

# Description

Helper func for calculating one-side overlap of SJ and junctions of gi

# Usage

```
.getStartEndOvl(gi, gr_chim_c, gr_sj_c, tol = 3)
```

# Arguments

gi GInteractions object

gr\_chim\_c GRanges object with chimeric junctions of gi

gr\_sj\_c GRanges with chimeric junctions

tol overlap tolerance

#### Value

Granges object with the left (A) region

.gv\_plotboxes 7

.gv\_plotboxes

Plots distributed boxes

# Description

Non-exported from Gviz, contains logic for calcualting boxes alignment on plot

# Usage

```
.gv_plotboxes(box, lwd, lty, alpha)
```

#### Value

pushes boxes to the viewport

 $.gv\_updatepars$ 

Set Gviz graphical parameters

# Description

Set Gviz graphical parameters

# Usage

```
.gv_updatepars(x, class)
```

### **Details**

Non-exported from Gviz. Uses the same procedure as in AnnotationTrack to set defaults.

# Value

set default values inside class upon calling withn class constructor

8 annotateGI

annota	т	e	l٦	П

Annotate RNA duplexes with features

# Description

Overlays RNA duplexes with GRanges annotation object.

# Usage

```
annotateGI(
   gi,
   anno_gr,
   keys = c("gene_name", "gene_type", "gene_id"),
   ambig_key = keys[1],
   save_ambig = TRUE,
   order_key = NULL,
   order_vec = NULL
)
```

# Arguments

gi	GInteraction object to annotate
anno_gr	GRanges object with the keys columns in the metadata
keys	vectors with the names of the features to use for annotation.
ambig_key	a which feature to use for recording the annotation ambiguity. Determines the values in $ambig_list.A$ and $ambig_list.B$ .
save_ambig	When RNA duplex overlaps multiple features, ambiguous annotation for a single key can be stored in $ambig_list.A$ and $ambig_list.B$ . ambig.A and 'ambig.B' fields will be added as the $0/1$ label
order_key	Experimental. In case RNA duplex overlaps multiple features, this key will be used to sort the overlapping features.
order_vec	Experimental. An ordered vector of values in order_key annotation feature, which sets the priority in case of feature overlap.

# **Details**

For each annotation feature in keys, i.e if keys=c(keyname1), then <keyname1>.A, <keyname1>.B annotation fields will be created, containing the names of overlapping features If no overlap is found for the feature, then filed will have NA

#### Value

GInteractions object with new fields

availableDisplayPars 9

### **Examples**

```
data("RNADuplexesSampleData")
annotateGI(gi = RNADuplexSampleDGs, anno_gr = SampleGeneAnnoGR)

# Prioritisation of the snRNA and lncRNA before mRNA if genes overlap
annotateGI(
    gi = RNADuplexSampleDGs, anno_gr = SampleGeneAnnoGR,
    keys = c("gene_id", "gene_name", "gene_type"),
    order_key = "gene_type",
    order_vec = c("snRNA", "lncRNA", "protein_coding"),
    save_ambig = TRUE
)
```

availableDisplayPars The default display parameters for a DuplexTrack object

#### **Description**

DuplexTrack inherits from [Gviz::Annotaiontrack()] and its Gviz parents. Most likely, user doesn't need all dioplay pars for the parents, so only parameters relevant to the DuplexTrack are returned by default.

#### Usage

```
availableDisplayPars(class)
```

#### **Arguments**

class

DuplexTrack track object This function allows user to display the default display parameters for the DuplexTrack class.

#### Value

list of the default display parameters.

```
library(InteractionSet)
anchor1 <- GRanges(
    seqnames = "chr1",
    ranges = IRanges(
        start = c(100, 600, 1100, 1600, 2100),
        end = c(200, 700, 1200, 1700, 2200)
    ),
    strand = "+"
)
anchor2 <- GRanges(
    seqnames = "chr1",
    ranges = IRanges(</pre>
```

```
start = c(300, 800, 1300, 1800, 2300),
    end = c(400, 900, 1400, 1900, 2400)
),
    strand = "+"
)

interactions <- GInteractions(anchor1, anchor2, mode = "strict")
gr_region <- range(anchor1, anchor2)
a <- DuplexTrack(interactions, gr_region = gr_region, stacking = "dense")
availableDisplayPars("DuplexTrack")
DuplexDiscovereR::availableDisplayPars(a)</pre>
```

calculateLigationPvalues

Calculate p-values and abundance fractions for RNA duplexes

# Description

Calculates p-values by applying Fisher test to each gene/transcript pair Uses BH correction, outputs duplex abundance relative to the per - gene/transcript count, and counts of other RNA duplexes formed by either or none gene/transcript in this pair.

# Usage

```
calculateLigationPvalues(gi, df_counts, id_col = "gene_id")
```

#### **Arguments**

gi	GInteraction object annotated with gene/transcript names
df_counts	data.frame A two- column dataframe with gene/transcript counts to. The first column should match the 'gene_id' feature in anno_gr. The second column is the respective count.
id_col	the prefix for gene/transcript metadata id fields in input gi. Two fields of $<$ id_col>.A and $<$ id.col>.B are expected. Otherwise throws error.

### **Details**

H0: RNA duplex not existing and reported due to the random ligation of fragments H1: RNA duplex is true and formed because of existing the RNA-RNA interaction

The probability of random ligation is modeled as (P(a, b)) given by the following equation: The probability P(a, b) is defined as:

$$P(a,b) \propto \begin{cases} 2 \cdot P(a) \cdot P(b) & \text{if } a:b \text{ is observed and } a \neq b \\ P(a) \cdot P(b) & \text{if } a:b \text{ is observed and } a=b \\ 0 & \text{else} \end{cases}$$

where The probability (P(a)) (same as for P(b) ) is calculated as:  $P(a) = \frac{N \, \text{reads}(a)}{\text{total N reads}}$ 

p-value calculated by comparing observed duplex abundance to the expected as the are under the curve distribution to the right of the observed. P(a, b) is normalized to sum up to one.

#### Value

GInteractions object with new fields

# **Examples**

```
data("RNADuplexesSampleData")
gi <- calculateLigationPvalues(RNADuplexSampleDGs, df_counts = RNADuplexesGeneCounts)
hist(gi$p.adj, breaks = 20)</pre>
```

classifyTwoArmChimeras

Wrapper for classification of the 2arm chimeric reads

### **Description**

Wraps two procedures for different types of classification for read alignment:

**overlap type** test if chimeric junction map to two non-overlapped regions or shorter than defined minimum distance

splice junction test if chimeric junction is also a splice junction

# Usage

```
classifyTwoArmChimeras(
   gi,
   min_junction_len = 4,
   junctions_gr,
   max_sj_shift = 4
)
```

# **Arguments**

#### **Details**

Calls detection of the chimeric junction type, annotates short junctions on same chromosome an strand as 'short'. Compares chimeric junctions with splice junctions. Adds results as the new metadata fields parallel to the input.

# Value

GInteractions object object of the same size with new columns:

```
splicejnc filled with 0 or 1
junction_type factor for the junction types
```

#### See Also

DuplexDiscovereR::getChimericJunctionTypes(), DuplexDiscovereR::getSpliceJunctionChimeras()

# **Examples**

clusterDuplexGroups

Cluster RNA duplexes in GInteractions object

#### **Description**

Main method to find duplex groups from the individual interactions

### Usage

```
clusterDuplexGroups(
   gi,
   graphdf = NULL,
   maxgap = 40,
   minoverlap = 10,
   id_column = "duplex_id",
   weight_column = "weight",
   fast_greedy = FALSE,
   decompose = FALSE,
   id_columns_grapdf = paste(id_column, c(1, 2), sep = "."),
   min_arm_ratio = 0.3,
   dump_graph = FALSE,
   dump_path = ""
)
```

clusterDuplexGroups 13

### **Arguments**

gi	GInteractions object
graphdf	Optional. Dataframe representing connection edges between entries in gi If not provided, graphdf is created inside the function
maxgap	For graph creation only. Max shift between arms starts and ends for pair of overlapping reads
minoverlap	For graph creation only. Minimum required overlap between either arm for pair of overlapping reads Other optional arguments, which are not relevant, unless user want to modify clustering weights or modify clustering in some other way
id_column	Optional. Column name in the GInteractions metadata, which was used to index temporary duplex groups, if they are present
weight_column	Optional. If graphdf is provided, field to use for weight overlaps
fast_greedy	Optional. Run the fast_greedy algorithm instead of Louvain. Can speed up calcualtion for the large graphs.
decompose	Decompose graph into separate sub-graphs before clustering.
id_columns_gra	pdf
	Column in the graph dataframe, which was used for index
min_arm_ratio	For graph creation only. Span-to-overlap ratio threshold. If smaller than this value, then edge is not drawn
dump_graph	For debug. Export the graph elements. not used
dump_path	For debug. PArt to export the graph elements. not used

#### **Details**

Accepts or creates the connections graphdf dataframe, creates graph with igraph package, uses community detection algoritm to call clusters. New field dg\_id is added to label the clusters (duplex groups). If no community is found for the read, dg\_id is NA

# Value

GInteractions object with new dg\_id column

collapseIdenticalReads

Collapses identical interactions

# Description

Two entries (reads) are considered identical if they share start, end, strand and score vales Identical entries are collapsed into the single one.

#### Usage

```
collapseIdenticalReads(gi)
```

# **Arguments**

gi

GInteractions(mode='strict') object with chromA, strandA, startA, endA, chromB, strandB, startB, endB, score columns Optionally cigar\_alnA, cigar\_alnB columns are also considered for collapsing 'read\_id' column used as the index in the initial objects. Created, if not exists

#### **Details**

Adds columns to the collapsed object duplex\_id (int) unique record id n\_reads (int) number of entries collapsed

#### Value

result\_list object with keys 'gi\_collapsed': New collapsed GInteraction object 'stats\_df': tibble with the mapping of the original entries to the new duplex\_id

```
# load data
data("RNADuplexesSmallGI")
res_collapse <- collapseIdenticalReads(SampleSmallGI)
gi_new <- res_collapse[["gi_collapsed"]]
# keeps the mapping of the colapsed object to new
read_stats_df <- res_collapse[["stats_df"]]</pre>
```

collapseSimilarChimeras

Call clustering multiple times to collapse similar reads into duplex groups

#### **Description**

Function calls clustering algorithm several times and collapses highly similar reads to the temporary duplex groups (DGs).

#### Usage

```
collapseSimilarChimeras(
   gi,
   read_stats_df,
   maxgap = 5,
   niter = 2,
   minoverlap = 10,
   min_nodes = 10
)
```

#### **Arguments**

gi	GInteractions	object

read\_stats\_df tibble with the mapping 'read\_id' and 'duplex\_id' fields 'read\_id' refers to the

unique read, 'duplex\_id' refers to the entry collapsed identical reads i.e two identical reads will will correspond to two unique read\_id and the single duplex\_id

with n\_reads=2

maxgap Maximum relative shift between the overlapping read arms

niter Number of times clustering will be called

minoverlap Minimum required overlap between either read arm

min\_nodes Minimum count of nodes to finish the interaction merging

#### Details

Calling this procedure before global read clustering substantially reduces time required for calling DGs. Collapsed duplex groups are aggregated only from the reads which are shifted by only a few nucleotides from each other. These DGs are temporary until full library clustering is called. To keep track of the mapping of the temprary DGs to the input, dedicated dataframe is returned. The 'duplex\_id' column will be added or updated as identifier for the temporary duplex group. The number of reads under single 'duplex\_id' is recorded in the 'n\_reads' fields

#### Value

a list with the following keys

gi\_updated GInteractions object with both collapsed duplex groups and not-collapsed unchanged reads

**stats\_df** tibble With the mapping from the unique read - with the infromation about time and memory reaquired for the function call

```
collapse_duplex_groups
```

Collapse the reads into the duplex groups after clustering

#### **Description**

Collapse each interaction in the input to the duplex group based on the pre-computed dg\_id

#### Usage

```
collapse_duplex_groups(
   gi,
   return_unclustered = FALSE,
   return_collapsed = TRUE,
   keep_meta = TRUE
)
```

### **Arguments**

gi GInteractions with the 'dg\_id' metadata field return\_unclustered

add unclustered reads to output

return\_collapsed

add duplex groups, which were created as temporary with n\_reads > 1 but was not clustered to the DG golabally. This parameter is used internally and should be kept default in most situations.

keep\_meta

whether to keep metadata, which only unclustered reads have, in case of a mixed output

# **Details**

'dg\_id' is used as the identifier for the duplex group Reads belonging to the same duplex group are collapsed into a single entry with start and end are set as min() and max() coordinate of the reads in within the duplex group. The 'score' column is averaged across the duplex group reads is calculated and put as the 'score' for the collapsed duplex group Behavior in case 'dg\_id' = NA: Option 'return\_unclustered' - whether unclustered reads with should be added to the output gi

**return\_unclustered == FALSE** Interaction is not returned in the output. Default.

col\_check\_rename 17

**return\_unclustered == TRUE** Interaction is returned in the output, output is mixed duplex groups and individual reads

Internally used argument #'

return\_collapsed == FALSE In case interaction already collapsed and n\_read > 1, interaction will not be returned as duplex group

**return\_collapsed == TRUE** In case interaction has n\_read > 1, interaction will be treated as duplex group

#### Value

GInteractions object with collapsed duplex groups

# **Examples**

```
# load example of clustered data
data("RNADuplexesSampleData")
# some reads assigned to DG, some are not
table(is.na(RNADuplexSampleGI$dg_id))
# Return only DGs
gicollapsed <- collapse_duplex_groups(RNADuplexSampleGI, return_unclustered = FALSE)</pre>
# Return DGs and unclustered reads as well
gimixed <- collapse_duplex_groups(RNADuplexSampleGI, return_unclustered = TRUE)</pre>
# load small sample GInteractions and process it manually
data("RNADuplexesSmallGI")
# First, collapse duplicated reads. This adds n_reads and duplex ids
ginodup <- collapseIdenticalReads(SampleSmallGI)$gi_collapsed</pre>
# Second, run clustering, get DG ids
ginodup <- clusterDuplexGroups(ginodup)</pre>
# Return all DGs result in n=3 DGS, one of them formed by
# identical duplicated alignments
collapse_duplex_groups(ginodup, return_collapsed = TRUE)
# Return DGs, but drop duplicated returns n=2 DGs
collapse_duplex_groups(ginodup, return_collapsed = FALSE)
```

col\_check\_rename

Check the column names and types in read dataframe

# **Description**

Function to check the correct column names and types in dataframe input. Tries to guess the column names, if colnames are not provided, but the types are correct

# Usage

```
col_check_rename(df, table_type = "STAR")
```

#### **Arguments**

```
df input table_type one in c("STAR","bedpe")
```

#### **Details**

- $\bullet \ \ Expected\ column\ names\ for\ bedge\ file\ c("chromA","startA", 'endA', "chromB", 'startB', 'endB', 'readname', 'endA', 'endA'$
- Expected colnames for STAR Chimeric junction input:
- For the 'old' chimeric detection scheme: c("chr\_donorA", "brkpt\_donorA", "strand\_donorA", "chr\_acceptorB",
- For the 'new' chimeric detection scheme c("chr\_donorA", "brkpt\_donorA", "strand\_donorA", "chr\_acceptorB", "brkpt\_acceptorB", "strand\_acceptorB", "junction\_type", "repeat\_left\_lenA", "repeat\_right\_lenB", "r "start\_alnA", "cigar\_alnA", "start\_alnB", "cigar\_alnB")

#### Value

dataframe with the properly formatted columns

```
compareMultipleInteractions
```

Compare multiple RNA-RNA interactions sets

# Description

Combines all interaction into single superset by clustering & collapsing. Then compares every input entry with the superset. Overlaps between superset and inputs are recorded in a table as 0/1

### Usage

```
compareMultipleInteractions(
  gi_samples_list,
  min_ratio = 0.3,
  minoverlap = 5,
  maxgap = 50,
  niter = 3,
  gi_superset = NULL,
  anno_gr = NULL
)
```

# Arguments

```
gi_samples_list
```

anmes list with the GInteractions entries list('sample1'=gi1,'sample2'='gi2)

min\_ratio

If the overlap-to-span ratio for either arm (A or B) for pair of chimeric reads is less than min\_arm\_ratio, then the total overlap for this pair is set to zero. Relevant to comparison of superset vs individual samples

minoverlap	Parameter for read clustering to create a superset. Minimum required overlap to for either arm (A or B) for pair of entries.
maxgap	Parameter for read clustering. Minimum required shift between start and end coordinates of arms for pair of overlapping entries If the shift is longer than max_gap for either arm, then total read overlap between those reads is zero.
niter	Internal parameter for debugging. Number of cluster& collapse iterations to find superset
gi_superset	Optional. Superset defining the space (all) of the interactions, against which inputs from the list will be compared.
anno_gr	Optional. Granges to annotate superset.

#### Value

dataframe recodding the overlaps between samples and supeset

```
# Create test set of RNA interactions
chrom <- "chr1"</pre>
start1 <- c(1, 11, 21, 31, 41, 51, 61, 71, 81, 91)
end1 <- start1 + 9
start2 <- c(101, 111, 121, 131, 141, 151, 161, 171, 181, 191)
end2 <- start2 + 9
anchor1 <- GRanges(seqnames = chrom, ranges = IRanges(start = start1, end = end1))</pre>
anchor2 <- GRanges(seqnames = chrom, ranges = IRanges(start = start2, end = end2))</pre>
interaction <- GInteractions(anchor1, anchor2)</pre>
# Ensure some overlaps
n <- length(interaction)</pre>
group_size <- ceiling(n / 2)</pre>
group_indices1 <- sort(sample(seq_len(n), group_size))</pre>
group_indices2 <- sort(sample(seq_len(n), group_size))</pre>
group_indices3 <- sort(sample(seq_len(n), group_size))</pre>
# Create separate GInteractions objects for each group
group1 <- interaction[group_indices1]</pre>
group2 <- interaction[group_indices2]</pre>
group3 <- interaction[group_indices3]</pre>
# format input and call comparison
a <- list("sample1" = group1, "sample2" = group2, "sample3" = group3)</pre>
res <- compareMultipleInteractions(a)</pre>
# comparison result
head(res$dt_upset)
# superset
res$gi_all
# dataframe for the Upset plot
res$dt_upset
```

20 convert\_gi\_to\_ranges

computeGISelfOverlaps Find overlaps between entries in GInteractions

#### **Description**

Utility function to find overlapping reads in the input and calculate overlap scores. Removes self-hits. Computes overlap/span ratios for each interaction arm. Sum of the scores is recorded in 'weight' field

# Usage

```
computeGISelfOverlaps(
   gi,
   id_column = "duplex_id",
   maxgap = 40,
   minoverlap = 10
)
```

# **Arguments**

gi input gi object

id\_column which use for using as ids for entries

maxgap parameter for call of InteractionSet::findOverlaps()
minoverlap parameter for call InteractionSet::findOverlaps()

#### Value

dataframe with indexes of pairwise overlapsin input and columns for span, overlap, ratios of either read arm

# **Examples**

```
data("RNADuplexesSmallGI")
computeGISelfOverlaps(SampleSmallGI)
```

```
convert_gi_to_ranges Convert GInteractions object to Granges
```

### Description

Creates the 'long' GRanges by stacking the A and B arms one 'on top' of the other. Adds id and group fields as indicators of original index and interaction arm (A- left arm, B- right arm)

#### Usage

```
convert_gi_to_ranges(gi)
```

dd\_get\_chimeric\_reads 21

# Arguments

```
gi GInteractions
```

#### Value

GRanges twice the length of the input

#### **Examples**

```
data("RNADuplexesSmallGI")
convert_gi_to_ranges(SampleSmallGI)
```

```
dd_get_chimeric_reads Accessor for chimeric_reads Slot
```

#### **Description**

Retrieves the value of the chimeric\_reads slot in a DuplexDiscovererResults object.

# Usage

```
dd_get_chimeric_reads(object)
## S4 method for signature 'DuplexDiscovererResults'
dd_get_chimeric_reads(object)
```

# Arguments

object

A DuplexDiscovererResults object.

#### Value

GInteractions object from the chimeric\_reads slot.

```
# load example input
data("RNADuplexesSmallGI")
data("RNADuplexesSampleData")
# run whole pipeline
result <- runDuplexDiscoverer(
   data = SampleSmallGI,
   junctions_gr = SampleSpliceJncGR,
   anno_gr = SampleGeneAnnoGR,
   sample_name = "run_example",
   lib_type = "SE",
   table_type = "STAR"
)
# access results</pre>
```

```
show(result)
gi_clusters <- dd_get_duplex_groups(result)
gi_reads <- dd_get_chimeric_reads(result)
df_reads <- dd_get_reads_classes(result)
dd_get_reads_classes(result)
dd_get_run_stats(result)</pre>
```

```
dd_get_chimeric_reads_stats
```

 $Accessor for \verb| chimeric_reads_stats| Slot$ 

#### **Description**

Retrieves the value of the chimeric\_reads\_stats slot in a DuplexDiscovererResults object.

#### Usage

```
dd_get_chimeric_reads_stats(object)
## S4 method for signature 'DuplexDiscovererResults'
dd_get_chimeric_reads_stats(object)
```

# Arguments

object

A DuplexDiscovererResults object.

### Value

tibble from the chimeric\_reads\_stats slot.

```
# load example input
data("RNADuplexesSmallGI")
data("RNADuplexesSampleData")
# run whole pipeline
result <- runDuplexDiscoverer(</pre>
   data = SampleSmallGI,
    junctions_gr = SampleSpliceJncGR,
    anno_gr = SampleGeneAnnoGR,
    sample_name = "run_example",
   lib_type = "SE",
    table_type = "STAR"
)
# access results
show(result)
gi_clusters <- dd_get_duplex_groups(result)</pre>
gi_reads <- dd_get_chimeric_reads(result)</pre>
df_reads <- dd_get_reads_classes(result)</pre>
dd_get_reads_classes(result)
dd_get_run_stats(result)
```

```
dd_get_duplex_groups Accessor for duplex_groups slot
```

# Description

Retrieves the value of the duplex\_groups slot in a DuplexDiscovererResults object.

# Usage

```
dd_get_duplex_groups(object)
## S4 method for signature 'DuplexDiscovererResults'
dd_get_duplex_groups(object)
```

# Arguments

object

A DuplexDiscovererResults object.

#### Value

GInteractions object from the duplex\_groups slot.

```
# load example input
data("RNADuplexesSmallGI")
data("RNADuplexesSampleData")
# run whole pipeline
result <- runDuplexDiscoverer(</pre>
   data = SampleSmallGI,
    junctions_gr = SampleSpliceJncGR,
    anno_gr = SampleGeneAnnoGR,
    sample_name = "run_example",
   lib_type = "SE",
    table_type = "STAR"
)
# access results
show(result)
gi_clusters <- dd_get_duplex_groups(result)</pre>
gi_reads <- dd_get_chimeric_reads(result)</pre>
df_reads <- dd_get_reads_classes(result)</pre>
dd_get_reads_classes(result)
dd_get_run_stats(result)
```

dd\_get\_reads\_classes

```
dd_get_reads_classes Accessor for reads_classes Slot
```

# Description

Retrieves the value of the reads\_classes slot in a DuplexDiscovererResults object.

# Usage

```
dd_get_reads_classes(object)
## S4 method for signature 'DuplexDiscovererResults'
dd_get_reads_classes(object)
```

# **Arguments**

object

A DuplexDiscovererResults object.

#### Value

tibble from the reads\_classes slot.

```
# load example input
data("RNADuplexesSmallGI")
data("RNADuplexesSampleData")
# run whole pipeline
result <- runDuplexDiscoverer(</pre>
   data = SampleSmallGI,
    junctions_gr = SampleSpliceJncGR,
    anno_gr = SampleGeneAnnoGR,
    sample_name = "run_example",
    lib\_type = "SE",
    table_type = "STAR"
)
# access results
show(result)
gi_clusters <- dd_get_duplex_groups(result)</pre>
gi_reads <- dd_get_chimeric_reads(result)</pre>
df_reads <- dd_get_reads_classes(result)</pre>
dd_get_reads_classes(result)
dd_get_run_stats(result)
```

dd\_get\_run\_stats 25

dd\_get\_run\_stats

Accessor for run\_stats Slot

# Description

Retrieves the value of the run\_stats slot in a DuplexDiscovererResults object.

# Usage

```
dd_get_run_stats(object)
## S4 method for signature 'DuplexDiscovererResults'
dd_get_run_stats(object)
```

# Arguments

object

A DuplexDiscovererResults object.

#### Value

tibble from the run\_stats slot.

```
# load example input
data("RNADuplexesSmallGI")
data("RNADuplexesSampleData")
# run whole pipeline
result <- runDuplexDiscoverer(</pre>
   data = SampleSmallGI,
    junctions_gr = SampleSpliceJncGR,
    anno_gr = SampleGeneAnnoGR,
    sample_name = "run_example",
    lib\_type = "SE",
    table_type = "STAR"
)
# access results
show(result)
gi_clusters <- dd_get_duplex_groups(result)</pre>
gi_reads <- dd_get_chimeric_reads(result)</pre>
df_reads <- dd_get_reads_classes(result)</pre>
dd_get_reads_classes(result)
dd_get_run_stats(result)
```

26 DuplexDiscovereR

drawGD, DuplexTrack-method

Draw method for DuplexTrack

# Description

Gviz::AnnotationTrack stacking algorithm is used to calculate vertical distribution of boxes for the interactions. Boxes coordinates are later imported for placing labels and arcs

### Usage

```
## S4 method for signature 'DuplexTrack'
drawGD(GdObject, minBase, maxBase, prepare = FALSE, subset = TRUE, ...)
```

#### Value

pushes boxes, arcs and labels to viewport

DuplexDiscovereR

Analysis of the data from RNA duplex probing experiments

#### **Description**

DuplexDiscovereR is a package for analysing data from RNA cross-linking and proximity ligation protocols such as SPLASH, PARIS, LIGR-seq and others, which provide information about intra-molecular RNA-RNA interactions through chimeric RNA-seq reads. Chimerically aligned fragments in these experiments correspond to the base-paired stretches (RNA duplexes) of RNA molecules . DuplexDiscovereR takes input in the form of chimericly or split -aligned reads, It implements procedures for alignment classification, filtering and efficient clustering of individual chimeric reads into duplex groups (DGs). Once DGs are found, RNA duplex formation and their hybridization energies are predicted. Additional metrics, such as p-values or mean DG alignment scores, can be calculated to rank and analyse the final set of RNA duplexes. Data from multiple experiments or replicates can be processed separately and further compared to check the reproducibility of the experimental method.

#### **Details**

**DuplexDiscovereR** 

#### Author(s)

Egor Semenchenko

### See Also

DuplexDiscovereR vignette

DuplexDiscovererResults-class

DuplexDiscovererResults

#### **Description**

A helper S4 class to store the results of the full DuplexDiscovereR analysis. This class contains the following output:

- duplex\_groups: clustered duplex groups.
- chimeric\_reads: individual two-regions chimeric reads. Contains both clustered and unclustered reads. Clustered reads are linked to the duplex groups though 'dg\_id' field in metadata
- reads\_classes: dataframe parallel to the input containing classification result and detected mapping type for each entry in the input
- chimeric\_reads\_stats: dataframe containing read type classification statistics
- run\_stats: data frame containing statistics about the time and memory used by the pipeline

# Usage

```
DuplexDiscovererResults(
  duplex_groups,
  chimeric_reads,
  reads_classes,
  chimeric_reads_stats,
  run_stats
)
```

#### **Arguments**

### **Details**

Each output type has a corresponding accessor:

```
dd_get_duplex_groups()dd_get_chimeric_reads()dd_get_reads_classes()dd_get_chimeric_reads_stats()dd_get_run_stats()
```

28 DuplexTrack

#### Value

A DuplexDiscovererResults object.

#### Slots

```
duplex_groups GInteractions object with duplex groups chimeric_reads GInteractions object with chimeric reads reads_classes tibble (tbl_df) with read classification data. chimeric_reads_stats tibble (tbl_df) read type statistics. run_stats tibble (tbl_df) runtime and memory info
```

#### See Also

```
dd_get_duplex_groups(),dd_get_chimeric_reads(),dd_get_reads_classes(),dd_get_chimeric_reads_stats()
,dd_get_run_stats()
```

#### **Examples**

```
# load example input
data("RNADuplexesSmallGI")
data("RNADuplexesSampleData")
# run whole pipeline
result <- runDuplexDiscoverer(</pre>
    data = SampleSmallGI,
    junctions_gr = SampleSpliceJncGR,
    anno_gr = SampleGeneAnnoGR,
    sample_name = "run_example",
    lib_type = "SE",
    table_type = "STAR"
)
# access results
show(result)
gi_clusters <- dd_get_duplex_groups(result)</pre>
gi_reads <- dd_get_chimeric_reads(result)</pre>
df_reads <- dd_get_reads_classes(result)</pre>
dd_get_reads_classes(result)
dd_get_run_stats(result)
```

DuplexTrack

class for the visualization of RNA duplexes

# Description

Inherits the Gviz::AnnotationTrack, plots interaction ranges as boxes. Arguments from Gviz::AnnotationTrack, as stacking which set boxes layout are accepted. Parent aesthetics for labels are overwritten with Display parameters of this class. Accepts GInteractions object to plot and GRanges to define plot region

DuplexTrack 29

Duplexes which can be displayed on the plot range are connected with arcs. Duplexes which are partially outside of the range are displayed without arcs. Labeles and appearance can be controlled with display parameters

#### **Arguments**

gi An GInteractions object gr\_region GRanges region for plotting

from Integer start coordinate of subset region. Used if gr\_region is not provided to Integer end coordinate of subset region. Used if gr\_region is not provided

chromosome Chromosome of subset region. Used if gr\_region is not provided

strand Used if gr\_region is not provided

fill.column used for fill. Default is "" (empty) and triggers IGV color pallete. **Display** parameters

arcs.color Character. Color of the arcs. Default is "black".

**arc.location** Character in c('inner','outer','midpoint'). Location of the arcs in X axis relative to range. Default is "inner"

**labels.v.offset.base** Numeric. Base vertical offset for the labels. Default is 0.2. Other offsets are added to it.

**labels.v.offset.trans** Numeric. Vertical offset for trans labels. Applied when one part of the duplex is outside of the plot. Recommended ranges are in -0.5 to 0.5 Default is 0.0.

**labels.h.offset.trans** Numeric. Horizontal offset for trans labels. Applied when one part of the duplex is outside of the plot Value is in nucleotide units. Default is 0.0.

**labels.v.offset.cis** Numeric. Vertical offset for cis labels. Recommended ranges are in -0.5 to 0.5 Default is 0.0. Default is 0.0.

**labels.h.offset.cis** Numeric. Horizontal offset for cis labels. Value is in nucleotide units. Default is 0.0.

labels.fontsize Numeric. Font size of the labels. Default is 18.

**label.cis.above** Logical. Whether the cis labels should be above. When set to FALSE, labels are plot for each box separately. Default is TRUE

**annotation.column1** Character. First annotation column to use for labels. Default is "group" and generated internally.

**annotation.column2** Character. Second annotation column to use for labels. Default is "" (empty).

**fill.column** Character. Column used for fill. Default is "" (empty) and triggers IGV color pallete.

**labels.color** Character. Color of the labels. Default is 'black'.

**labels.align** Character. Alignment of the labels. Default is 'center'. Possible values are in c('left','right','center)

arcConstrain Numeric. Minimum gap distance between arms of the interaction to draw arcs

30 DuplexTrack

```
library(InteractionSet)
library(Gviz)
# generate input
anchor1 <- GRanges(</pre>
    seqnames = "chr1",
    ranges = IRanges(
        start = c(100, 600, 1100, 1600, 2100, 150, 400),
        end = c(200, 700, 1200, 1700, 2200, 250, 500)
    ),
    strand = "+"
)
anchor2 <- GRanges(</pre>
    segnames = "chr1",
    ranges = IRanges(
        start = c(300, 800, 1300, 1800, 2300, 1500, 1700),
        end = c(400, 900, 1400, 1900, 2400, 1600, 1800)
    ),
    strand = "+"
)
interactions <- GInteractions(anchor1, anchor2, mode = "strict")</pre>
# define plotting range
gr_region <- range(anchor1, anchor2)</pre>
interactions$anno_A <- sample(LETTERS, length(interactions))</pre>
interactions$anno_B <- interactions$anno_A</pre>
a <- DuplexTrack(interactions, gr_region = gr_region, stacking = "dense")
plotTracks(a, stacking = "dense")
plotTracks(a, stacking = "squish", annotation.column1 = "anno_A")
# add interactions which are not fully in plot range: outside the range or on different chromosome()
# one left (A) interaction arm outside of the plot, other on different chromosome
new_anchor1 <- GRanges(</pre>
    seqnames = c("chr1", "chr2"),
    ranges = IRanges(
        start = c(10, 600),
        end = c(90, 700)
    ),
    strand = "+"
)
new_anchor2 <- GRanges(</pre>
    seqnames = c("chr1", "chr1"),
    ranges = IRanges(
        start = c(1500, 1000),
        end = c(1600, 1200)
    ),
    strand = "+"
)
new_interactions <- GInteractions(new_anchor1, new_anchor2)</pre>
new_interactions$anno_A <- c("A.out", "A.out_chr")</pre>
```

```
new_interactions$anno_B <- c("B.in", "B.in")
all_interactions <- c(interactions, new_interactions)
b <- DuplexDiscovereR::DuplexTrack(all_interactions,
    gr_region = gr_region,
    annotation.column1 = "anno_A",
    annotation.column2 = "anno_B"
)
plotTracks(b)
# to customize plot, one can call, to see options
DuplexDiscovereR::availableDisplayPars(b)</pre>
```

getChimericJunctionTypes

Classify chimeric junctions of two-arm reads into types

### **Description**

Chimeric reads which can be represented ans two-arm interactions can be divided into several categories based on the distance between the chimeric fragments and existence of the overlap between these fragments.

# Usage

```
getChimericJunctionTypes(gi, normal_gap_threshold = 10)
```

### **Arguments**

```
gi GInteractions object
normal_gap_threshold
minimum allowed distance between chimeric arms
```

#### **Details**

Takes GInteractions object and classifies junctions into following categories

```
2arm normal chimeric read
2arm_short normal chimeric read with junction < normal_gap_threshold
self_ovl arms overlap
antisense_ovl arms overlap on the opposite strand</pre>
```

#### Value

```
gi object of the same size with the 'junction_type' field added
```

32 getRNAHybrids

### **Examples**

```
data("RNADuplexesSampleData")
preproc_df <- runDuplexDiscoPreproc(RNADuplexesRawBed, table_type = "bedpe")
preproc_gi <- makeGiFromDf(preproc_df)
preproc_gi <- getChimericJunctionTypes(preproc_gi)
table(preproc_gi$junction_type)</pre>
```

getRNAHybrids

Run prediciton of RNA hybridization

### **Description**

Calls RNAduplex from ViennaRNA to find base-pairs for every entry in the input, throws a message and system warning if it is not installed

### Usage

```
getRNAHybrids(gi, fafile)
```

# **Arguments**

```
gi Ginteraction with pairs of regions fafile path to the .fasta file with genome
```

#### Value

object parallel to input with added energy GC content, dot-format base-pairings and lenghts of RNA hybrids will return the input, if RNAhybrids cannot be run

```
sequence <- paste0(</pre>
   "CGUAGCAUCGUAGCUAGCUAGCUAUGCGAUU"
)
# Save the sequence to a temp fasta file
fasta_file <- tempfile(fileext = ".fa")</pre>
chrom <- "test_chrA"</pre>
writeLines(c(">test_chrA", sequence), con = fasta_file)
# Create the GInteraction object
# Define start and end positions for the base-pairing regions
regions <- data.frame(</pre>
   start1 = c(1, 11, 21, 31, 41),
   end1 = c(10, 20, 30, 40, 50),
   start2 = c(91, 81, 71, 61, 51),
   end2 = c(100, 90, 80, 70, 60)
)
```

```
# GRanges objects for the anchors
anchor1 <- GRanges(seqnames = chrom, ranges = IRanges(start = regions$start1, end = regions$end1))
anchor2 <- GRanges(seqnames = chrom, ranges = IRanges(start = regions$start2, end = regions$end2))
interaction <- GInteractions(anchor1, anchor2)
# predict hybrids
# In case ViennaRNA is installed
## Not run:
gi_with_hybrids <- getRNAHybrids(interaction, fasta_file)
## End(Not run)</pre>
```

getSpliceJunctionChimeras

Identify chimeric junctions coinciding with the splice junctions

# **Description**

Marks interactions which starts/ends within specified shift from the known splice junctions.

### Usage

```
getSpliceJunctionChimeras(
   gi,
   sj_gr,
   sj_tolerance = 20,
   sj_tolerance_strict = 10
)
```

# Arguments

gi	GInteractions object		
sj_gr	Granges object with the splice junctions data		
sj_tolerance	maximum shift between either donor and acceptor splice sites and corresponding chimreic junction coordinates to count chimeric junction as splice junction		
sj_tolerance_strict			

maximum shift between either donor and acceptor splice sites irrespective of the particular splice junction. If both chimeric junction start and end correspond to donor or acceptor of any known junction, it is marked as splice junction. Used to catch novel combinations of known 3' and 5' sites

#### Value

gi object with added 'splicejnc' and field Additionally 'splicejnc\_donor' 'splicejnc\_acceptor' fields are added

34 get\_arm\_b

# **Examples**

```
data("RNADuplexesSampleData")
gi <- getSpliceJunctionChimeras(RNADuplexSampleGI, SampleSpliceJncGR)
table(gi$splicejnc)
table(gi$splicejnc_acceptor, gi$splicejnc_donor)</pre>
```

get\_arm\_a

Get left arm of GInteraction

# Description

Get left arm of GInteraction

#### Usage

```
get_arm_a(gi)
```

# Arguments

gi

GInteractions object

#### Value

Granges object with the left (A) region

get\_arm\_b

Get right arm of GInteraction

# Description

Get right arm of GInteraction

# Usage

```
get_arm_b(gi)
```

# **Arguments**

gi

GInteractions object

# Value

Granges object with the right (B) region

get\_char\_count\_cigar 35

```
get_char_count_cigar Count the length of the key type in CIGAR string
```

# **Description**

Takes CIGAR operands i.e M,N,S and sums the associated blocks length It is vectorized. i.e supports vector with CIGAR strings

#### Usage

```
get_char_count_cigar(strings, s)
```

# **Arguments**

strings CIGAR string vector s CIGAR operands

#### Value

vector with length values

# **Examples**

```
# From a vector
get_char_count_cigar(c("4S18M22S", "25S26M"), "S")
get_char_count_cigar(c("18M22S", "20M20S"), "M")
```

# Description

Returns chimeric junction defined as range distance between the end and the start of the first and second range respectively

# Usage

```
get_chimeric_junctions_onestrand(gi_intra)
```

#### **Details**

If the pair of interacting ranges is not on the same strand and chromosome, returns error

#### Value

Granges object with the

36 makeDfFromGi

# **Description**

Get colnames for expected data types

### Usage

```
get_colnames_and_types_for_input(nameset)
```

# **Arguments**

nameset

name of the table

#### Value

character vector

makeDfFromGi

Convert GInteractions to tibble

# **Description**

Converts GInteractions to tibble, preserves metadata

### Usage

```
makeDfFromGi(gi)
```

# **Arguments**

gi

**GInteracttions** 

# **Details**

```
Following naming conventions is used for region coordinates: c('chromA', 'startA', 'endA', 'strandA', 'chromB', 'startB', 'endB', 'strandB')
```

#### Value

tibble preserving metadata columns

#### See Also

```
makeGiFromDf()
```

makeGiFromDf 37

## **Examples**

```
data(RNADuplexesSmallGI)
converted_to_df <- makeDfFromGi(SampleSmallGI)
converted_to_gi <- makeGiFromDf(converted_to_df)</pre>
```

makeGiFromDf

Convert Dataframe to GInteractions

## **Description**

Converts dataframe-like object to the GInteractions.

# Usage

```
makeGiFromDf(df)
```

#### **Arguments**

df

dataframe-like object. Should be convertable to tibble::tibble()

#### **Details**

arms will be consistent between different objects of same reference Following columns are looked up in input dataframe to parse region coordinates: c("chromA','startA','endA','strandA',"chromB",'startB','endB','strandB') GInteractions (mode='strict') is enforced, to ensure that the order of the regions Extra columns are stored as metadata fields

## Value

GInteractions(mode='strict')

#### See Also

```
makeDfFromGi()
```

# **Examples**

```
# load example GInteractions
data(RNADuplexesSmallGI)

converted_to_df <- makeDfFromGi(SampleSmallGI)
converted_to_gi <- makeGiFromDf(converted_to_df)</pre>
```

```
preproc_chim_junction_out_pe
```

Processing of of the STAR PE Chimeric.junction.out

#### **Description**

Calculates alignment coordinates and returns reads with categories

## Usage

```
preproc_chim_junction_out_pe(dt, keep_all_columns = FALSE)
```

## **Arguments**

dt Chimeric.out.junction with the correct column names keep\_all\_columns

• TRUE or FALSE. Keep CIGAR strings and junction coordinate columns

#### **Details**

#

multimap multi-mapped read

multigap more than one junction (more than two 'N' in CIGAR string)

**bad junction** Artifacts. I.e alignments for both arms are continious, but with 'backward' chimeric junction was wrongly put

#### Value

tibble with annotated reads

```
preproc_chim_junction_out_se
```

Processing of of the STAR SE Chimeric.junction.out

## **Description**

Calculates alignment coordinates and returns reads with categories

## Usage

```
preproc_chim_junction_out_se(dt, keep_all_columns = FALSE)
```

preproc\_generic 39

## Arguments

```
dt Chimeric.out.junction with the correct column names keep_all_columns
```

• TRUE or FALSE. Keep CIGAR strings and junction coordinate columns

#### **Details**

#'

multimap multi-mapped read

multigap more than one junction (more than two 'N' in CIGAR string)

**bad junction** Artifacts. I.e alignments for both arms are continious, but with 'backward' chimeric junction was wrongly put

#### Value

tibble with annotated reads

#### See Also

```
col_check_rename()
```

preproc\_generic

Preprocess .bedpe input

# Description

Searches for the multi-mapped and bad reads (overlapping arms) Adds 'multimap', 'bad\_junction' columns filled with 0 or 1 and 'multigap' = 0 for consistency with other pre-processing methods

## Usage

```
preproc_generic(dt, keep_all_columns = TRUE)
```

## **Arguments**

```
dt dataframe with reads aligned to strictly two loci keep_all_columns
```

keeep columns apart form the required from .bedpe format

## Value

pre-processed dataframe

40 refresh\_gi

preproc\_generic\_gi

Preprocess GInteractions input

# Description

Searches for the multi-mapped reads (overlapping arms) Adds 'multimap', 'bad\_junction' columns filled with 0/1 and 'multigap' = 0 for consistency with other pre-processing methods.

## Usage

```
preproc_generic_gi(gi_raw, keep_all_columns = TRUE)
```

## **Arguments**

```
gi_raw GInteractions with inpit RNA interactions keep_all_columns
```

keep columns apart from those required by .bedpe format

#### Value

pre-processed dataframe

refresh\_gi

Refresh the GInteractions object

## Description

Sub-setting the GInteractions object does not reduce its ranges container For some applications, to save memory, we can safely reduce the size of the object by re-creating it. Also, it can be used to ensure the 'strict' mode of the regions in ranges

## Usage

```
refresh_gi(gi)
```

## Arguments

gi

GInteractions object

## Value

GInteractions object with new ranges attribute

RNADuplexesGeneCounts Gene counts on human chromosome 22, embryonic stem cells

## **Description**

File generated by mapping with STAR using --quantMode GeneCounts see system.file("extdata/scripts", "DD\_data\_generation.R", package = "DuplexDiscovereR") for details on the pre-processing and sub-setting the

# Usage

data(RNADuplexesSampleData)

#### **Format**

An object of class spec\_tbl\_df (inherits from tbl\_df, tbl, data.frame) with 1445 rows and 2 columns.

#### Value

tibble with columns of Chimeric.junction.out

#### **Source**

SequenceReadArcive

 ${\tt RNADuplexesRawBed}$ 

Chimeric reads of SPLASH converted to .bedpe fromat

## **Description**

A Chimeric.out.Junction file with a subset of chr 22 Chimeric reads detected by SPLASH protocol in Human embryonic stem cells.

# Usage

data(RNADuplexesSampleData)

#### **Format**

An object of class spec\_tbl\_df (inherits from tbl\_df, tbl, data.frame) with 2040 rows and 10 columns.

### Value

tibble with columns of bedpe format

#### Source

SequenceReadArcive Reads were aligned with STAR and filtered to contain only reads which could be represented as 2-arm chimeric alignments. Converted to the <a href="bedpe">bedpe</a> format see system.file("extdata/scripts", "DD\_data\_generation.R", package = "DuplexDiscovereR") for details on the pre-processing and sub-setting the data

RNADuplexesRawChimSTAR

Chimeric reads of SPLASH

## Description

A Chimeric out. Junction file with a subset of chr 22 Chimeric reads detected by SPLASH protocol in Human embryonic stem cells.

## Usage

data(RNADuplexesSampleData)

#### **Format**

An object of class tbl\_df (inherits from tbl, data.frame) with 5000 rows and 21 columns.

## Value

tibble with columns of Chimeric.junction.out

## Source

SequenceReadArcive Reads were aligned with STAR see system.file("extdata/scripts", "DD\_data\_generation.R", package = "DuplexDiscovereR") for details on the pre-processing and sub-setting the data

RNADuplexSampleClustReads

RNA duplex reads of SPLASH, clustered and assigned to duplex groups

## **Description**

GInteractions read-level object containing processed reads, annotated with duplex group ids, read types gene names and p-values

## Usage

data(RNADuplexesSampleData)

#### **Format**

An object of class StrictGInteractions of length 2090.

#### Value

GInteractions with

- n\_reads\_dg : number of reads in the duplex group (DG)
- duplex\_id : temporary id for RNA duplexes which could be found before clustering (duplicated or shifted by couple of nt )
- dg\_id :id of the duplex group
- score: median alignment score in duplex group
- other columns inherited from the STAR Chimeric.out.Junction

#### Source

SequenceReadArcive Reads were aligned with STAR and duplex groups were identified see system.file("extdata/script"DD\_data\_generation.R", package = "DuplexDiscovereR") for details on the data generation proceedure.

RNADuplexSampleDGs

RNA duplex reads of SPLASH, clustered and collapsed to duplex groups

# Description

GInteractions duplex group -level object containing detected duplex groups, annotated with duplex group ids, gene\_names and p-values

## Usage

data(RNADuplexesSampleData)

## Format

An object of class StrictGInteractions of length 79.

#### Value

GInteractions with

- n\_reads : number of reads in the duplex group (DG)
- dg\_id :id of the duplex group
- p\_val: BH adjusted p-value of testing to reject hypothesis of DG arising from random ligation
- score: median alignment score in duplex group
- other columns with . A and . B annotating to which genes either arm of the DG maps

#### Source

SequenceReadArcive Reads were aligned with STAR and duplex groups were identified see system.file("extdata/script"DD\_data\_generation.R", package = "DuplexDiscovereR") for details on the data generation procedure.

RNADuplexSampleGI

RNA duplex reads of SPLASH derived from chimeric alignments

# Description

GInteractions read-level object containing two-arm chimeric reads extracted from mapping output and which can be represented in the GInteraction object

#### Usage

data(RNADuplexesSampleData)

#### **Format**

An object of class StrictGInteractions of length 2090.

#### Details

see system.file("extdata/scripts", "DD\_data\_generation.R", package = "DuplexDiscovereR") for details on the data generation proceedure.

#### Value

GInteractions with

- readname: read name
- map\_type : type of the mapped read (2arm by design of pre-filtering)
- junction\_type : if read junction is too short, or it not a 'true' ligated reads because of the junction coincides with splice junction
- cigar\_aln\* columns inherited from the STAR Chimeric.out.Junction output

#### Source

SequenceReadArcive

runDuplexDiscoPreproc Run pre-processing of chimeric reads input

#### Description

Imports dataframe with reads (.bedpe or Chimeric.out.junction) or GInteractions object. Checks column names or tries to quess them if not provided. Adds necessary annotation depending on the input type, For STAR input, calculates length of the alignments and marks unique 2-arm alignments. For the .bedpe or GInteractions input, all entries are already represented as reads with two different aligned parts (2-arm), so only check for unique readname is performed.

# Usage

```
runDuplexDiscoPreproc(
  data,
  table_type,
  library_type = "SE",
  keep_metadata = TRUE,
  return_gi = FALSE,
  min_arm_len = 15
)
```

### **Arguments**

data	Either dataframe-like object: $Chimeric.out.junction\ from\ STAR$ or $.bedpe$ - formatted or GInteractions object from $InteractionSet$ package
table_type	in c("STAR", "bedpe") for Chimeric.out.Junction or generic input
library_type	c("SE", "PE") for pair- or single- end input
keep_metadata	${\tt c(TRUE,FALSE)}$ Whether extra fields like CIGAR strings and junction coordinates should be kept
return_gi	if the return object should be GInteractions
min_arm_len	minimum allowed length of the alignment arm. Read will be dropped if either arm is shorter

## **Details**

If not existed, adds fields required for the downstream steps: 'readname', 'map\_type', 'score', 'n\_reads'. 'map\_type' field determines the type of the chimeric read:

multimap multi-mapped read

multigap more than one junction (more than two 'N' in CIGAR string)

**bad junction** Artifacts or possibly unaccounted types. I.e alignments for both arms are continuous, but with 'backward' chimeric junction was wrongly introduced in the mapping

46 runDuplexDiscoverer

#### Value

tibble with new metadata fields OR GInteractions if return\_gi is set to TRUE

## **Examples**

runDuplexDiscoverer

Executes all steps of DuplexDiscovereR pipeline

#### **Description**

Generates GInteractions object with duplex groups from the STAR Chimeric.out.junction or bedpe file. Classifies reads, annotates reads by overlap with the gene or transcript features, calculates p-values and hybridization energies. Additionally, returns mappings from duplex groupd back to genes.

## Usage

```
runDuplexDiscoverer(
  data,
  table_type = "",
  junctions_gr = NULL,
  anno_gr = NULL,
  anno_gr_keys = c("gene_id", "gene_name", "gene_type"),
  fafile = NULL,
  df_counts = NULL,
  sample_name = "sample",
  lib_type = "SE",
 min_junction_len = 5,
 max_gap = 50,
 min_arm_ratio = 0.1,
 min_overlap = 10,
 max_sj_shift = 10,
  gap_collapse_similar = 2,
  collapse_n_inter = 5,
  trim_alignments = FALSE,
  trim_length = 40,
 min_arm_len = 9,
  compute_p_values = TRUE
)
```

47 runDuplexDiscoverer

#### **Arguments**

data dataframe-like object with the split reads. Output of Chimeric.out.junction or

dataframe with fileds defined by bedpe format: c("chromA", "startA", 'endA', "chromB", 'startB', 'endB', 'rea

... ) Alternatively, GInteractions object

one in c("STAR", "bedpe") Defines the type of the input dataframe. ignored if table\_type

input data is GInteractions

**GRanges** object with the splice junction coordinates junctions\_gr

GRanges object to use for the annotation of the interactions. Optional anno\_gr

c() vector with names of metadata fields in anno gr which will be used for the anno\_gr\_keys

annotation. Argument passed to annotateGI() function. The c('gene\_id','gene\_name','gene\_type')

columns in anno gr are used by default.

fafile path to the genome .fasta file. Used to calculate hybridization energy with

RNADuplex. Sequence names should correspond to the sequences from which

the mapping index was created. Optional

df\_counts A two- column dataframe with counts. Counts are used for p-value calculation.

The first column should match the 'gene\_id' feature in anno\_gr. The second

column is the respective count. Optional

sample\_name A name of the sample, used for assembling the analysis statistics dataframe

one in c('SE','PE'). Type of the sequencing library. Default is 'SE' lib\_type

min\_junction\_len

a minimum allowed distance between chimeric arms for the read input. Reads with the junction closer than min\_junction\_len are annotated as '2arm' shot'

and not clustered to duplex groups

Parameter for read clustering. Minimum required shift between start and end max\_gap

> coordinates of arms for pair of overlapping chimeric reads. If the shift is longer than max\_gap for either arm, then total read overlap between those reads is zero.

Parameter for read clustering. If the overlap-to-span ratio for either arm (A or min\_arm\_ratio

B) for pair of chimeric reads is less than min\_arm\_ratio, then the total overlap

for this pair is set to zero.

Parameter for read clustering. Minimum required overlap to for either arm (A min\_overlap

or B) for pair of chimeric reads.

Maximum shift between either donor and acceptor splice sites and chimeric max\_sj\_shift

junction coordinates to count chimeric junction as splice junction

gap\_collapse\_similar

Parameter for read clustering (iterative step). Analogous to the max\_gap, but applied collapse\_n\_inter times during the iterative merging step. Reduce this to 1 or 2 to lower RAM usage for clustering the library with many similar

reads.

collapse\_n\_inter

Parameter for read clustering (iterative step). Number of iterations to repeat step of collapsing of the highly similar chimeric reads. Increasing this from i.e 0 to 5 reduces clustering time and memory for the libraries with many overlapping

reads.

48 runDuplexDiscoverer

trim\_alignments

TRUE or FALSE. Whether to trim arms alignments to 'trim\_length' nucleotide

around chimeric junction

trim\_length target size of trimmed alignment

min\_arm\_len minimum allowed length of the alignment arm. Read will be dropped if either

arm is shorter

compute\_p\_values

TRUE or FALSE. whether to calcualte random ligation test

#### **Details**

This is a main function to do the initial discovery of the RNA duplexes after the chimeric read mapping. It wraps following procedures:

- Classifies the input reads by the mapping type. Keeps 2-arm chimeric reads for downstream analysis
- Compares 2arm duplex reads against provided splice junctions
- Classifies 2arm duplexes into spurious self-overlapping, splice junction categoris
- Performs clustering of the remaining reads into duplex groups
  - Collapses identically mapped reads
  - Collapses closely located reads, almost identical reads
  - Finds duplex groups throughout whole data set
- Annotates duplex groups with genomic features if annotation is provided
- Calculates p-values if gene counts and annotation are provided
- Calculates hybridization energies if path to the .fasta file is provided

#### Value

a list with the following keys

duplex\_groups GInteractions object with chimeric reads clustered duplex groups

chimeric\_reads GInteractions object with non-collapsed chimeric reads

reads\_classes tbl\_df dataframe parallel to the the input dataframe, annotated with read categories and duplex groups

chimeric\_reads\_stats tbl\_df dataframe containing read type classification statistics

run\_stats tbl\_df dataframe with the time and memory info about the run

#### See Also

DuplexDiscovererResults()

SampleGeneAnnoGR 49

#### **Examples**

```
library(DuplexDiscovereR)
# load data
data("RNADuplexesSampleData")
result <- runDuplexDiscoverer(</pre>
   data = RNADuplexesRawChimSTAR,
    junctions_gr = SampleSpliceJncGR,
    anno_gr = SampleGeneAnnoGR,
    df_counts = RNADuplexesGeneCounts,
    sample_name = "test clustering",
    fafile = NULL,
   collapse_n_inter = 3,
   lib_type = "SE",
    table_type = "STAR"
# see results object
print(result)
# duplex groups
dd_get_duplex_groups(result)
# individual chimeric reads
dd_get_chimeric_reads(result)
# counts of detected read types
dd_get_chimeric_reads_stats(result)
```

SampleGeneAnnoGR

Gene coordinates on human chromosome 22

## Description

Granges containing gene coordinates of human chromosome 22 obtained from GENCODEv44 annotaion

## Usage

```
data(RNADuplexesSampleData)
```

#### **Format**

An object of class GRanges of length 1445.

## **Details**

```
see system.file("extdata/scripts", "DD_data_generation.R", package = "DuplexDiscovereR")
for details
```

### Value

```
statdatd GENCODE gtf fields
```

50 SampleSpliceJncGR

#### **Source**

**GENCODEv44** 

SampleSmallGI

RNA duplex reads of SPLASH derived from chimeric alignments

# Description

GInteractions object containing two-arm chimeric reads extracted from mapping output and which can be represented in the GInteraction object and subset to chr22: 23877144-45562960 '\*,

## Usage

data(RNADuplexesSmallGI)

#### **Format**

An object of class StrictGInteractions of length 14.

#### **Details**

see system.file("extdata/scripts", "DD\_data\_generation.R", package = "DuplexDiscovereR") for details on the data generation procedure.

## **Source**

SequenceReadArcive

SampleSpliceJncGR

Gene coordinates on human chromosome 22

## Description

Granges containing coordinates of splice junctions human chromosome 22 obtained from GEN-CODEv44 annotaion

# Usage

data(RNADuplexesSampleData)

### **Format**

An object of class GRanges of length 8465.

## **Details**

```
see system.file("extdata/scripts", "DD_data_generation.R", package = "DuplexDiscovereR")
for details
```

#### Value

statdatd GENCODE gtf fields

#### **Source**

**GENCODEv44** 

 $\verb|show,DuplexDiscovererResults-method|\\$ 

Show Method for DuplexDiscovererResults class

# Description

This method provides a summary of the DuplexDiscovererResults object. It prints chimeric\_reads\_stats followed by the run\_stats.

#### Usage

```
## S4 method for signature 'DuplexDiscovererResults'
show(object)
```

## **Arguments**

object

A DuplexDiscovererResults object.

#### Value

None. Prints a formatted summary.

show, DuplexTrack-method

Show method for DuplexTrack

## **Description**

Show method for DuplexTrack

#### **Usage**

```
## S4 method for signature 'DuplexTrack'
show(object)
```

52 subset\_gi

## **Arguments**

object DuplexTrack.

#### Value

class representation

# **Examples**

```
library(InteractionSet)
anchor1 <- GRanges(</pre>
   seqnames = "chr1"
   ranges = IRanges(
        start = c(100, 600, 1100, 1600, 2100),
        end = c(200, 700, 1200, 1700, 2200)
   ),
    strand = "+"
)
anchor2 <- GRanges(</pre>
    segnames = "chr1",
   ranges = IRanges(
        start = c(300, 800, 1300, 1800, 2300),
        end = c(400, 900, 1400, 1900, 2400)
   ),
    strand = "+"
)
interactions <- GInteractions(anchor1, anchor2, mode = "strict")</pre>
gr_region <- range(anchor1, anchor2)</pre>
a <- DuplexTrack(interactions, gr_region = gr_region, stacking = "dense")</pre>
show(a)
```

subset\_gi

Subset the GInteractions object to single interaction

# Description

Sub-setting the GInteractions object does not reduce its ranges container This function selects Ginteraction by index and reduces the ranges

## Usage

```
subset_gi(gi, k)
```

# Arguments

gi

GInteractions object

trimAroundJunction 53

## Value

GInteractions with the range atribure reduced to single interaction

trimAroundJunction

Extract regions around chimeric junction

# Description

Trim alignements to contain only 'extract len' nucleotides adajcent to the chimeric junction

## Usage

```
trimAroundJunction(dt, extract_len = 30)
```

## **Arguments**

```
dt table with the
extract_len the length of the sequence to trim
```

## **Details**

In case of the long alignemtns, it may be necessary trim chimeric alignments to identify RNA duplex. If 'extract\_len' is longer than the read alignemnt length, then no trimming is performed

## Value

dataframe with the trimmed alignments

## **Examples**

```
data("RNADuplexesSampleData")
dt_preproc <- runDuplexDiscoPreproc(RNADuplexesRawChimSTAR,
        table_type = "STAR", library_type = "SE"
)
trimAroundJunction(dt_preproc, 40)</pre>
```

54 writeGiToSAMfile

writeGiToSAMfile

Write reads to sam file

## **Description**

Writes interactions to the sam file for visualization in extrnal browsers. Takes input as GInteractions object containing reads or duplex groups.

## Usage

```
writeGiToSAMfile(
  gi_coords,
  file_out,
  distance_chim_junction = 10000,
  read_name_column = "readname",
  id_column = "dg_id",
  genome = "",
  sample_name = "noname_sample"
)
```

#### **Arguments**

gi\_coords input Ginteraction object file\_out path to write output file distance\_chim\_junction

maximum distance between input duplex groups/reads, which will be represented as the single-line in .sam file. Junction will be output as N- gap. For the interactions with longer distances, chimeric junction will be represented as MR:Z:i tag

read\_name\_column

character field, pointing out to read names. Read names are generated automatically if not provided.

id\_column character name of the field containing integer duplex group ids. NA are replaced

with zeros

genome character. Genome version. Required for the retrieval of sequence lengths for

sam file header- SQ and SN tags. For convenience, hg38 and hg19 chromosome lengths will be assigned automatically. If the value is not in c('hg38','hg19'), seqlengths will be looked for be in attribute in seqlengths() of regions(gi\_coords)

sample\_name name to use in RG SAM tag in header

#### Value

no object is returned

writeGiToSAMfile 55

## **Examples**

```
# Load test data
data("RNADuplexesSampleData")
# if the input is read-based, it should have integer duplex group ids
# here, we have 2090 reads
length(RNADuplexSampleGI)
\# among them 300 reads does not belong to any DG
\# missing ids will be converted to 0
table(is.na(RNADuplexSampleGI$dg_id))
tmpf <- tempfile(".sam")</pre>
writeGiToSAMfile(
   gi_coords = RNADuplexSampleGI,
    id_column = "dg_id",
    file_out = tmpf,
   distance_chim_junction = 1e5,
   genome = "hg38"
)
```

# **Index**

* datasets	.gv_updatepars,7	
RNADuplexesGeneCounts, 41		
RNADuplexesRawBed, 41	annotateGI, $8$	
RNADuplexesRawChimSTAR, 42	availableDisplayPars, $9$	
RNADuplexSampleClustReads, 42		
RNADuplexSampleDGs, 43	calculateLigationPvalues, $10$	
RNADuplexSampleGI, 44	classifyTwoArmChimeras, 11	
SampleGeneAnnoGR, 49	clusterDuplexGroups, 12	
SampleSmallGI, 50	col_check_rename, 17	
SampleSpliceJncGR, 50	<pre>col_check_rename(), 39</pre>	
* interal	collapse_duplex_groups, 16	
<pre>preproc_chim_junction_out_se, 38</pre>	collapseIdenticalReads, 14	
* internal	collapseSimilarChimeras, 15	
.DGIdToDuplexId, 6	${\sf compareMultipleInteractions}, 18$	
.annotateCisTrans, 5	computeGISelfOverlaps, 20	
.compute_clusters_comp, 5	convert_gi_to_ranges, 20	
.getStartEndOvl, 6	dd mae alaimania maada 01	
.gv_plotboxes, 7	dd_get_chimeric_reads, 21	
.gv_updatepars,7	dd_get_chimeric_reads(), 27, 28	
col_check_rename, 17	dd_get_chimeric_reads, DuplexDiscovererResults-method	
drawGD, DuplexTrack-method, 26	<pre>(dd_get_chimeric_reads), 21 dd_get_chimeric_reads_stats, 22</pre>	
get_arm_a, 34	dd_get_chimeric_reads_stats(), 27, 28	
get_arm_b, 34		
<pre>get_chimeric_junctions_onestrand,</pre>	<pre>dd_get_chimeric_reads_stats,DuplexDiscovererResults-method</pre>	
35	(dd_get_chilleric_reads_stats), 22	
<pre>get_colnames_and_types_for_input,</pre>	dd_get_duplex_groups, 23	
36	dd_get_duplex_groups(), 27, 28	
<pre>preproc_chim_junction_out_pe, 38</pre>	dd_get_duplex_groups,DuplexDiscovererResults-method	
preproc_generic, 39	(dd_get_duplex_groups), 23	
preproc_generic_gi, 40	dd_get_reads_classes, 24	
refresh_gi, 40	dd_get_reads_classes(), 27, 28	
subset_gi, 52	dd_get_reads_classes,DuplexDiscovererResults-method	
.DGIdToDuplexId, 6	(dd_get_reads_classes), 24	
.addDGidsForTmpDGs,3	dd_get_run_stats, 25	
.addGeneCounts, 4	dd_get_run_stats(), 27, 28	
.annotateCisTrans, 5	dd_get_run_stats,DuplexDiscovererResults-method	
.compute_clusters_comp, 5	(dd_get_run_stats), 25	
.getStartEndOvl, 6	drawGD, DuplexTrack-method, 26	
gy plothoxes 7	DunleyDiscovereR 26	

INDEX 57

```
DuplexDiscovereR-package
                                                trimAroundJunction, 53
        (DuplexDiscovereR), 26
{\tt DuplexDiscovereR::getChimericJunctionTypes(),writeGiToSAMfile,54}
        DuplexDiscovereR::getSpliceJunctionChimeras(),
        12
DuplexDiscovererResults
        (DuplexDiscovererResults-class),
        27
DuplexDiscovererResults(), 48
DuplexDiscovererResults-class, 27
DuplexTrack, 28
get_arm_a, 34
get_arm_b, 34
get_char_count_cigar, 35
get_chimeric_junctions_onestrand, 35
{\tt get\_colnames\_and\_types\_for\_input, 36}
getChimericJunctionTypes, 31
getRNAHybrids, 32
getSpliceJunctionChimeras, 33
InteractionSet::findOverlaps(), 20
makeDfFromGi, 36
makeDfFromGi(), 37
makeGiFromDf, 37
makeGiFromDf(), 36
preproc_chim_junction_out_pe, 38
preproc_chim_junction_out_se, 38
preproc_generic, 39
preproc_generic_gi, 40
refresh_gi, 40
RNADuplexesGeneCounts, 41
RNADuplexesRawBed, 41
RNADuplexesRawChimSTAR, 42
RNADuplexSampleClustReads, 42
RNADuplexSampleDGs, 43
RNADuplexSampleGI, 44
runDuplexDiscoPreproc, 45
runDuplexDiscoverer, 46
SampleGeneAnnoGR, 49
SampleSmallGI, 50
SampleSpliceJncGR, 50
show, DuplexDiscovererResults-method,
        51
show, DuplexTrack-method, 51
subset_gi, 52
```