

# Package: DuplexDiscoverer (via r-universe)

October 31, 2024

**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** 1.1.0

**Roxygen** list(markdown = TRUE)

**RoxygenNote** 7.3.1

**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

**Depends** R (>= 4.4), InteractionSet

**LazyData** false

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

**VignetteBuilder** knitr  
**VignetteEngine** knitr  
**Config/testthat/edition** 3  
**Repository** <https://bioc.r-universe.dev>  
**RemoteUrl** <https://github.com/bioc/DuplexDiscovererR>  
**RemoteRef** HEAD  
**RemoteSha** 37477e10cc8137ff130f14ec40e1c062fafa0a93

## Contents

.addDGidsForTmpDGs . . . . .	3
.addGeneCounts . . . . .	4
annotateGI . . . . .	4
availableDisplayPars . . . . .	5
classifyTwoArmChimeras . . . . .	6
clusterDuplexGroups . . . . .	7
collapseIdenticalReads . . . . .	9
collapseSimilarChimeras . . . . .	10
collapse_duplex_groups . . . . .	11
compareMultipleInteractions . . . . .	12
computeGISelfOverlaps . . . . .	14
convert_gi_to_ranges . . . . .	15
dd_get_chimeric_reads . . . . .	15
dd_get_chimeric_reads_stats . . . . .	16
dd_get_duplex_groups . . . . .	17
dd_get_reads_classes . . . . .	18
dd_get_run_stats . . . . .	19
DuplexDiscovererR . . . . .	20
DuplexDiscovererResults-class . . . . .	20
DuplexTrack . . . . .	22
getChimericJunctionTypes . . . . .	25
getRNAHybrids . . . . .	26
getSpliceJunctionChimeras . . . . .	27
get_char_count_cigar . . . . .	28
makeDfFromGi . . . . .	28
makeGiFromDf . . . . .	29
preproc_chim_junction_out_se . . . . .	30
RNADuplexesGeneCounts . . . . .	30
RNADuplexesRawBed . . . . .	31
RNADuplexesRawChimSTAR . . . . .	32
RNADuplexSampleClustReads . . . . .	32
RNADuplexSampleDGs . . . . .	33
RNADuplexSampleGI . . . . .	34
runDuplexDiscoPreproc . . . . .	34
runDuplexDiscoverer . . . . .	36
SampleGeneAnnoGR . . . . .	38

<code>.addDGidsForTmpDGs</code>	3
SampleSmallGI . . . . .	39
SampleSpliceIncGR . . . . .	40
show,DuplexDiscovererResults-method . . . . .	40
show,DuplexTrack-method . . . . .	41
writeGiToSAMfile . . . . .	42
<b>Index</b>	<b>44</b>

---

<code>.addDGidsForTmpDGs</code>	<i>Helper function to add ids to the duplex groups missed during global clustering</i>
---------------------------------	--

---

**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`

**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`

---

<code>.addGeneCounts</code>	<i>Helper function to add count data to metadata of GInteractions</i>
-----------------------------	---

---

**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**

<code>gi</code>	GInteractions
<code>df_counts</code>	dataframe with read counts
<code>id_col</code>	key to use in merge

**Value**

GInteractions with added counts

---

<code>annotateGI</code>	<i>Annotate RNA duplexes with features</i>
-------------------------	--

---

**Description**

Overlays RNA duplexes with GRanges annotation object.

**Usage**

```
annotateGI(
  gi,
  anno_gr,
  keys = c("gene_name", "gene_type", "gene_id"),
  save_ambig = TRUE
)
```

**Arguments**

<code>gi</code>	GInteraction object to annotate
<code>anno_gr</code>	GRanges object with the keys columns in the metadata
<code>keys</code>	names of the features to use for annotation.
<code>save_ambig</code>	in case RNA duplex overlaps multiple features of the first key, mark the existence of ambiguous annotation in the fields <code>ambig.A</code> and <code>ambig.B</code> . Fields <code>ambig_list.A</code> and <code>ambig_list.B</code> will be store the list of overlapping features. Only the first filed from keys is checked for possible annotation ambiguities.

**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 field will have NA

**Value**

GInteractions object with new fields

**Examples**

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

---

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.

**Examples**

```
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(
  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")
DuplexDiscovererR::availableDisplayPars(a)

```

---

```
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

<code>gi</code>	GInteractions object
<code>min_junction_len</code>	minimum allowed distance between two chimeric arms
<code>junctions_gr</code>	Granges object with the splice junctions coordinates
<code>max_sj_shift</code>	maximum shift between either donor and acceptor splice sites and corresponding chimreic junction coordinates to count chimeric junction as splice junction

## 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 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**

```
data("RNADuplexesSampleData")
head(RNADuplexSampleGI)
# remove all metadata
mcols(RNADuplexSampleGI) <- NULL
gi <- classifyTwoArmChimeras(RNADuplexSampleGI,
  min_junction_len = 5,
  junctions_gr = SampleSpliceJncGR, max_sj_shift = 10
)
table(gi$splicejnc)
table(gi$junction_type)
```

---

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 = ""
)
```

**Arguments**

<code>gi</code>	GInteractions object
<code>graphdf</code>	Optional. Dataframe representing connection edges between entries in <code>gi</code> . If not provided, <code>graphdf</code> is created inside the function
<code>maxgap</code>	For graph creation only. Max shift between arms starts and ends for pair of overlapping reads
<code>minoverlap</code>	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
<code>id_column</code>	Optional. Column name in the GInteractions metadata, which was used to index temporary duplex groups, if they are present
<code>weight_column</code>	Optional. If <code>graphdf</code> is provided, field to use for weight overlaps
<code>fast_greedy</code>	Optional. Run the <code>fast_greedy</code> algorithm instead of Louvain. Can speed up calculation for the large graphs.
<code>decompose</code>	Decompose graph into separate sub-graphs before clustering.
<code>id_columns_graphdf</code>	Column in the graph dataframe, which was used for index
<code>min_arm_ratio</code>	For graph creation only. Span-to-overlap ratio threshold. If smaller than this value, then edge is not drawn
<code>dump_graph</code>	For debug. Export the graph elements. not used
<code>dump_path</code>	For debug. Path to export the graph elements. not used

**Details**

Accepts or creates the connections `graphdf` dataframe, creates graph with `igraph` package, uses community detection algorithm 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

**Examples**

```
data("RNADuplexesSampleData")
# run preprocessing and filtering
preproc_df <- runDuplexDiscoPreproc(RNADuplexesRawBed, table_type = "bedpe")
preproc_gi <- makeGiFromDf(preproc_df)
preproc_gi <- classifyTwoArmChimeras(preproc_gi,
  min_junction_len = 5,
  junctions_gr = SampleSpliceJncGR, max_sj_shift = 10
)
# collapse duplicates
gi <- collapseIdenticalReads(preproc_gi)$gi
# run global clustering
gi <- clusterDuplexGroups(gi)
# check dg_ids
table(is.na(gi$dg_id))
```



---

`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`

## Examples

```
# 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"]]
```

---

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
)
```

## 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 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

## 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 temporary 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 the information about time and memory required 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 globally. 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.

**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)
# Return DGs and unclustered reads as well
gimixed <- collapse_duplex_groups(RNADuplexSampleGI, return_unclustered = TRUE)

# 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
# Second, run clustering, get DG ids
ginodup <- clusterDuplexGroups(ginodup)
# 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)
```

---

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**

<code>gi_samples_list</code>	names list with the GInteractions entries list('sample1'=gi1,'sample2'=gi2)
<code>min_ratio</code>	If the overlap-to-span ratio for either arm (A or B) for pair of chimeric reads is less than <code>min_arm_ratio</code> , then the total overlap for this pair is set to zero. Relevant to comparison of superset vs individual samples
<code>minoverlap</code>	Parameter for read clustering to create a superset. Minimum required overlap to for either arm (A or B) for pair of entries.
<code>maxgap</code>	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 <code>max_gap</code> for either arm, then total read overlap between those reads is zero.
<code>niter</code>	Internal parameter for debugging. Number of cluster& collapse iterations to find superset
<code>gi_superset</code>	Optional. Superset defining the space (all) of the interactions, against which inputs from the list will be compared.
<code>anno_gr</code>	Optional. Granges to annotate superset.

**Value**

dataframe recodding the overlaps between samples and supeset

**Examples**

```
# Create test set of RNA interactions
chrom <- "chr1"
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))
anchor2 <- GRanges(seqnames = chrom, ranges = IRanges(start = start2, end = end2))

interaction <- GInteractions(anchor1, anchor2)

# Ensure some overlaps
n <- length(interaction)
group_size <- ceiling(n / 2)
group_indices1 <- sort(sample(seq_len(n), group_size))
group_indices2 <- sort(sample(seq_len(n), group_size))
group_indices3 <- sort(sample(seq_len(n), group_size))

# Create separate GInteractions objects for each group
group1 <- interaction[group_indices1]
group2 <- interaction[group_indices2]
group3 <- interaction[group_indices3]

# format input and call comparison
a <- list("sample1" = group1, "sample2" = group2, "sample3" = group3)
```

```
res <- compareMultipleInteractions(a)
# comparison result
head(res$dt_upset)
# superset
res$gi_all
# dataframe for the Upset plot
res$dt_upset
```

---

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	column which use for using as ids for entries
maxgap	parameter for call of <a href="#">InteractionSet::findOverlaps()</a>
minoverlap	parameter for call <a href="#">InteractionSet::findOverlaps()</a>

## Value

dataframe with indexes of pairwise overlaps in 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)
```

**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.

## Examples

```
# 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
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)
```

---

dd\_get\_chimeric\_reads\_stats

*Accessor for 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.



## Examples

```
# 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
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)
```

---

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.

## Examples

```
# 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
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)

```

---

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.

### Examples

```

# 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
  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)
```

---

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.

### Examples

```
# 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
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)
```

---

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 chimerically 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 through 'dg\_id' field in metadata

- reads\_classes: dataframe parallel to the 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

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

### 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\(\)](#)

### 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(
  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)
gi_reads <- dd_get_chimeric_reads(result)
df_reads <- dd_get_reads_classes(result)
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

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. Labels and appearance can be controlled with display parameters

**Arguments**

<code>gi</code>	An <code>GInteractions</code> object
<code>gr_region</code>	<code>GRanges</code> region for plotting
<code>from</code>	Integer start coordinate of subset region. Used if <code>gr_region</code> is not provided
<code>to</code>	Integer end coordinate of subset region. Used if <code>gr_region</code> is not provided
<code>chromosome</code>	Chromosome of subset region. Used if <code>gr_region</code> is not provided

strand	Used if gr_region is not provided
fill.column	used for fill. Default is "" (empty) and triggers IGV color palette. <b>Display parameters</b>
	<b>arcs.color</b> Character. Color of the arcs. Default is "black".
	<b>arc.location</b> Character in c('inner','outer','midpoint'). Location of the arcs in X axis relative to range. Default is "inner"
	<b>labels.v.offset.base</b> Numeric. Base vertical offset for the labels. Default is 0.2. Other offsets are added to it.
	<b>labels.v.offset.trans</b> 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.
	<b>labels.h.offset.trans</b> 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.
	<b>labels.v.offset.cis</b> Numeric. Vertical offset for cis labels. Recommended ranges are in -0.5 to 0.5 Default is 0.0. Default is 0.0.
	<b>labels.h.offset.cis</b> Numeric. Horizontal offset for cis labels. Value is in nucleotide units. Default is 0.0.
	<b>labels.fontsize</b> Numeric. Font size of the labels. Default is 18.
	<b>label.cis.above</b> Logical. Whether the cis labels should be above. When set to FALSE, labels are plot for each box separately. Default is TRUE
	<b>annotation.column1</b> Character. First annotation column to use for labels. Default is "group" and generated internally.
	<b>annotation.column2</b> Character. Second annotation column to use for labels. Default is "" (empty).
	<b>fill.column</b> Character. Column used for fill. Default is "" (empty) and triggers IGV color palette.
	<b>labels.color</b> Character. Color of the labels. Default is 'black'.
	<b>labels.align</b> Character. Alignment of the labels. Default is 'center'. Possible values are in c('left','right','center')
	<b>arcConstrain</b> Numeric. Minimum gap distance between arms of the interaction to draw arcs

## Examples

```

library(InteractionSet)
library(Gviz)
# generate input
anchor1 <- GRanges(
  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(

```

```

    seqnames = "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")
# define plotting range
gr_region <- range(anchor1, anchor2)
interactions$anno_A <- sample(LETTERS, length(interactions))
interactions$anno_B <- interactions$anno_A
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(
  seqnames = c("chr1", "chr2"),
  ranges = IRanges(
    start = c(10, 600),
    end = c(90, 700)
  ),
  strand = "+"
)
new_anchor2 <- GRanges(
  seqnames = c("chr1", "chr1"),
  ranges = IRanges(
    start = c(1500, 1000),
    end = c(1600, 1200)
  ),
  strand = "+"
)

new_interactions <- GInteractions(new_anchor1, new_anchor2)
new_interactions$anno_A <- c("A.out", "A.out_chr")
new_interactions$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)

```



---

`getChimericJunctionTypes`*Classify chimeric junctions of two-arm reads into types*

---

### Description

Chimeric reads which can be represented as 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

### Value

gi object of the same size with the 'junction\_type' field added

### 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)
```



```
gi_with_hybrids <- getRNAHybrids(interaction, fasta_file)

## End(Not run)
```

---

```
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

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

### Examples

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

---

`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

<code>strings</code>	CIGAR string vector
<code>s</code>	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")
```

---

`makeDfFromGi` *Convert GInteractions to tibble*

---

### Description

Converts GInteractions to tibble, preserves metadata

### Usage

```
makeDfFromGi(gi)
```

### Arguments

<code>gi</code>	GInteractions
-----------------	---------------

### 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\(\)](#)

**Examples**

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

---

makeGiFromDf	<i>Convert Dataframe to GInteractions</i>
--------------	---

---

**Description**

Converts dataframe-like object to the GInteractions.

**Usage**

```
makeGiFromDf(df)
```

**Arguments**

df                    dataframe-like object. Should be convertible 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)
```

---

```
preproc_chim_junction_out_se
```

*Processing 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)
```

### 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 continuous, but with 'backward' chimeric junction was wrongly put

### Value

tibble with annotated reads

### See Also

[col\\_check\\_rename\(\)](#)

---

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](#)

---

RNADuplexesRawBed	<i>Chimeric reads of SPLASH converted to .bedpe format</i>
-------------------	--

---

**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 `bedpe` 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**

[SequenceReadArchive](#) Reads were aligned with STAR see `system.file("extdata/scripts", "DD_data_generation.R", package = "DuplexDiscovererR")` 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**

**SequenceReadArchive** 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.

---

RNADuplexSampleDGs	<i>RNA duplex reads of SPLASH, clustered and collapsed to duplex groups</i>
--------------------	---

---

**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**

**SequenceReadArchive** 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	<i>RNA duplex reads of SPLASH derived from chimeric alignments</i>
-------------------	--

---

### 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 procedure.

### 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

[SequenceReadArchive](#)

---

runDuplexDiscoPreproc	<i>Run pre-processing of chimeric reads input</i>
-----------------------	---

---

### Description

Imports dataframe with reads (*.bedpe* or *Chimeric.out.junction*) or GInteractions object. Checks column names or tries to guess 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
)
```

**Arguments**

data	Either dataframe-like object: <i>Chimeric.out.junction</i> from <i>STAR</i> or <i>.bedpe</i> - formatted or GInteractions object from <b>InteractionSet</b> 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	c(TRUE, FALSE) Whether extra fields like CIGAR strings and junction coordinates should be kept
return_gi	if the return object should be GInteractions

**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

**Value**

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

**Examples**

```
# load data
data(RNADuplexesSampleData)
# with bedpe input
preproc_reads <- runDuplexDiscoPreproc(RNADuplexesRawBed, table_type = "bedpe")
# with STAR input
preproc_reads_star <- runDuplexDiscoPreproc(RNADuplexesRawChimSTAR,
  table_type = "STAR",
  keep_metadata = FALSE
)
```

---

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 group back to genes.

## Usage

```
runDuplexDiscoverer(
  data,
  table_type = "",
  junctions_gr = NULL,
  anno_gr = NULL,
  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
)
```

## Arguments

data	dataframe-like object with the split reads. Output of Chimeric.out.junction or dataframe with fields defined by bedpe format: c("chromA","startA","endA","chromB","startB","endB",...) Alternatively, GInteractions object
table_type	one in c("STAR","bedpe") Defines the type of the input dataframe. ignored if input data is GInteractions
junctions_gr	<b>GRanges</b> object with the splice junction coordinates
anno_gr	<b>GRanges</b> object to use for the annotation of the interactions. The c('gene_id','gene_name','gene_types') columns in anno_gr are used by default. Optional
fafile	path to the genome .fasta file. Used to calculate hybridization energy with <i>RNADuplex</i> . Sequence names should correspond to the sequences from which the mapping index was created. Optional
df_counts	A two- column dataframe with counts to use 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
lib_type	one in c('SE','PE'). Type of the sequencing library. Default is 'SE'
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
max_gap	Parameter for read clustering. Minimum required shift between start and end 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.
min_arm_ratio	Parameter for read clustering. 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.
min_overlap	Parameter for read clustering. Minimum required overlap to for either arm (A or B) for pair of chimeric reads.
max_sj_shift	Maximum shift between either donor and acceptor splice sites and chimeric 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.

## 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 DuplexDiscovererResults with the following output

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\(\)](#)

**Examples**

```
library(DuplexDiscoverer)
# load data
data("RNADuplexesSampleData")
result <- runDuplexDiscoverer(
  data = RNADuplexesRawChimSTAR,
  junctions_gr = SampleSpliceJncGR,
  anno_gr = SampleGeneAnnoGR,
  df_counts = RNADuplexesGeneCounts,
  sample_name = "test clustering",
  fofile = 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 tyoes
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

**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 GENCODEv44 annotation

**Usage**

```
data(RNADuplexesSampleData)
```

**Format**

An object of class GRanges of length 8465.

**Details**

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

**Value**

statdatd GENCODE gtf fields

**Source**

[GENCODEv44](#)

---

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)
```

**Arguments**

object            DuplexTrack.

**Value**

class representation

**Examples**

```
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(  
    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")  
show(a)
```

---

writeGiToSAMfile      *Write reads to sam file*

---

### Description

Writes interactions to the sam file for visualization in external 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

**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")
writeGiToSAMfile(
  gi_coords = RNADuplexSampleGI,
  id_column = "dg_id",
  file_out = tmpf,
  distance_chim_junction = 1e5,
  genome = "hg38"
)
```

# Index

- \* **datasets**
  - [RNADuplexesGeneCounts](#), [30](#)
  - [RNADuplexesRawBed](#), [31](#)
  - [RNADuplexesRawChimSTAR](#), [32](#)
  - [RNADuplexSampleClustReads](#), [32](#)
  - [RNADuplexSampleDGs](#), [33](#)
  - [RNADuplexSampleGI](#), [34](#)
  - [SampleGeneAnnoGR](#), [38](#)
  - [SampleSmallGI](#), [39](#)
  - [SampleSpliceJncGR](#), [40](#)
- \* **internal**
  - [preproc\\_chim\\_junction\\_out\\_se](#), [30](#)
  - [.addDGidsForTmpDGs](#), [3](#)
  - [.addGeneCounts](#), [4](#)
- [annotateGI](#), [4](#)
- [availableDisplayPars](#), [5](#)
  
- [classifyTwoArmChimeras](#), [6](#)
- [clusterDuplexGroups](#), [7](#)
- [col\\_check\\_rename\(\)](#), [30](#)
- [collapse\\_duplex\\_groups](#), [11](#)
- [collapseIdenticalReads](#), [9](#)
- [collapseSimilarChimeras](#), [10](#)
- [compareMultipleInteractions](#), [12](#)
- [computeGISelfOverlaps](#), [14](#)
- [convert\\_gi\\_to\\_ranges](#), [15](#)
  
- [dd\\_get\\_chimeric\\_reads](#), [15](#)
- [dd\\_get\\_chimeric\\_reads\(\)](#), [21](#), [22](#)
- [dd\\_get\\_chimeric\\_reads, DuplexDiscovererResults-method](#)  
([dd\\_get\\_chimeric\\_reads](#)), [15](#)
- [dd\\_get\\_chimeric\\_reads\\_stats](#), [16](#)
- [dd\\_get\\_chimeric\\_reads\\_stats\(\)](#), [21](#), [22](#)
- [dd\\_get\\_chimeric\\_reads\\_stats, DuplexDiscovererResults-method](#)  
([dd\\_get\\_chimeric\\_reads\\_stats](#)), [16](#)
- [dd\\_get\\_duplex\\_groups](#), [17](#)
- [dd\\_get\\_duplex\\_groups\(\)](#), [21](#), [22](#)
  
- [dd\\_get\\_duplex\\_groups, DuplexDiscovererResults-method](#)  
([dd\\_get\\_duplex\\_groups](#)), [17](#)
- [dd\\_get\\_reads\\_classes](#), [18](#)
- [dd\\_get\\_reads\\_classes\(\)](#), [21](#), [22](#)
- [dd\\_get\\_reads\\_classes, DuplexDiscovererResults-method](#)  
([dd\\_get\\_reads\\_classes](#)), [18](#)
- [dd\\_get\\_run\\_stats](#), [19](#)
- [dd\\_get\\_run\\_stats\(\)](#), [21](#), [22](#)
- [dd\\_get\\_run\\_stats, DuplexDiscovererResults-method](#)  
([dd\\_get\\_run\\_stats](#)), [19](#)
- [DuplexDiscovererR](#), [20](#)
- [DuplexDiscovererR-package](#)  
([DuplexDiscovererR](#)), [20](#)
- [DuplexDiscovererR::getChimericJunctionTypes\(\)](#),  
[DuplexDiscovererR::getSpliceJunctionChimeras\(\)](#),  
[7](#)
- [DuplexDiscovererResults](#)  
([DuplexDiscovererResults-class](#)),  
[20](#)
- [DuplexDiscovererResults\(\)](#), [38](#)
- [DuplexDiscovererResults-class](#), [20](#)
- [DuplexTrack](#), [22](#)
  
- [get\\_char\\_count\\_cigar](#), [28](#)
- [getChimericJunctionTypes](#), [25](#)
- [getRNAHybrids](#), [26](#)
- [getSpliceJunctionChimeras](#), [27](#)
  
- [InteractionSet::findOverlaps\(\)](#), [14](#)
  
- [makeDfFromGi](#), [28](#)
- [makeDfFromGi\(\)](#), [29](#)
- [makeGiFromDf](#), [29](#)
- [makeGiFromDf\(\)](#), [29](#)
  
- [preproc\\_chim\\_junction\\_out\\_se](#), [30](#)
- [RNADuplexesGeneCounts](#), [30](#)
- [RNADuplexesRawBed](#), [31](#)
- [RNADuplexesRawChimSTAR](#), [32](#)
- [RNADuplexSampleClustReads](#), [32](#)

RNADuplexSampleDGs, [33](#)  
RNADuplexSampleGI, [34](#)  
runDuplexDiscoPreproc, [34](#)  
runDuplexDiscoverer, [36](#)

SampleGeneAnnoGR, [38](#)  
SampleSmallGI, [39](#)  
SampleSpliceJncGR, [40](#)  
show, DuplexDiscovererResults-method,  
[40](#)  
show, DuplexTrack-method, [41](#)

writeGiToSAMfile, [42](#)