

# Package ‘annotatr’

March 29, 2021

**Title** Annotation of Genomic Regions to Genomic Annotations

**Version** 1.16.0

**Date** 2020-07-07

**Description** Given a set of genomic sites/regions (e.g. ChIP-seq peaks, CpGs, differentially methylated CpGs or regions, SNPs, etc.) it is often of interest to investigate the intersecting genomic annotations. Such annotations include those relating to gene models (promoters, 5'UTRs, exons, introns, and 3'UTRs), CpGs (CpG islands, CpG shores, CpG shelves), or regulatory sequences such as enhancers. The annotatr package provides an easy way to summarize and visualize the intersection of genomic sites/regions with genomic annotations.

**Depends** R (>= 3.4.0)

**Imports** AnnotationDbi, AnnotationHub, dplyr, GenomicFeatures, GenomicRanges, GenomeInfoDb (>= 1.10.3), ggplot2, IRanges, methods, readr, regioneR, reshape2, rtracklayer, S4Vectors (>= 0.23.10), stats, utils

**Suggests** BiocStyle, devtools, knitr, org.Dm.eg.db, org.Gg.eg.db, org.Hs.eg.db, org.Mm.eg.db, org.Rn.eg.db, rmarkdown, roxygen2, testthat, TxDb.Dmelanogaster.UCSC.dm3.ensGene, TxDb.Dmelanogaster.UCSC.dm6.ensGene, TxDb.Ggallus.UCSC.galGal5.refGene, TxDb.Hsapiens.UCSC.hg19.knownGene, TxDb.Hsapiens.UCSC.hg38.knownGene, TxDb.Mmusculus.UCSC.mm9.knownGene, TxDb.Mmusculus.UCSC.mm10.knownGene, TxDb.Rnorvegicus.UCSC.rn4.ensGene, TxDb.Rnorvegicus.UCSC.rn5.refGene, TxDb.Rnorvegicus.UCSC.rn6.refGene

**VignetteBuilder** knitr

**BugReports** <https://www.github.com/rcavalcante/annotatr/issues>

**License** GPL-3

**LazyData** true

**RoxygenNote** 7.1.0

**biocViews** Software, Annotation, GenomeAnnotation, FunctionalGenomics, Visualization

**git\_url** <https://git.bioconductor.org/packages/annotatr>

**git\_branch** RELEASE\_3\_12

**git\_last\_commit** e5b2b37

**git\_last\_commit\_date** 2020-10-27

**Date/Publication** 2021-03-29

**Author** Raymond G. Cavalcante [aut, cre],  
Maureen A. Sartor [ths]

**Maintainer** Raymond G. Cavalcante <rcavalca@umich.edu>

## R topics documented:

annotate_regions . . . . .	3
annotations . . . . .	4
annotatr . . . . .	4
annotatr_cache . . . . .	5
build_ah_annots . . . . .	5
build_annotations . . . . .	6
build_cpg_annots . . . . .	7
build_enhancer_annots . . . . .	8
build_gene_annots . . . . .	8
build_hmm_annots . . . . .	9
build_lncrna_annots . . . . .	9
builtin_annotations . . . . .	10
builtin_genomes . . . . .	10
check_annotations . . . . .	11
expand_annotations . . . . .	11
get_cellline_from_code . . . . .	12
get_cellline_from_shortcut . . . . .	12
get_orgdb_name . . . . .	13
get_txdb_name . . . . .	13
plot_annotation . . . . .	14
plot_categorical . . . . .	15
plot_coannotations . . . . .	17
plot_numerical . . . . .	19
plot_numerical_coannotations . . . . .	21
randomize_regions . . . . .	23
read_annotations . . . . .	24
read_regions . . . . .	25
reformat_hmm_codes . . . . .	26
subset_order_tbl . . . . .	26
summarize_annotations . . . . .	27
summarize_categorical . . . . .	28
summarize_numerical . . . . .	29
tidy_annotations . . . . .	30

annotate\_regions      *A function to intersect user region data with annotation data*

**Description**

Annotate genomic regions to selected genomic annotations while preserving the data associated with the genomic regions.

**Usage**

```

annotate_regions(
  regions,
  annotations,
  minoverlap = 1L,
  ignore.strand = TRUE,
  quiet = FALSE
)

```

**Arguments**

- regions            The GRanges object returned by read\_regions().
- annotations      A character vector of annotations to build. Valid annotation codes are listed with builtin\_annotations(). The "basicgenes" shortcut builds the following regions: 1-5Kb upstream of TSSs, promoters, 5UTRs, exons, introns, and 3UTRs. The "cpgs" shortcut builds the following regions: CpG islands, shores, shelves, and interCGI regions. NOTE: Shortcuts need to be appended by the genome, e.g. hg19\_basicgenes. Custom annotations whose names are of the form [genome]\_custom\_[name] should also be included. Custom annotations should be read in and converted to GRanges with read\_annotations(). They can be for a supported\_genome(), or for an unsupported genome.
- minoverlap        A scalar, positive integer, indicating the minimum required overlap of regions with annotations.
- ignore.strand     Logical indicating whether strandedness should be respected in findOverlaps(). Default FALSE.
- quiet             Print progress messages (FALSE) or not (TRUE).

**Value**

A GRanges where the granges are from the regions, and the mcols include the mcols from the regions and a column with the annotation GRanges.

**Examples**

```

r_file = system.file('extdata', 'test_read_multiple_data_nohead.bed', package='annotatr')
extraCols = c(pval = 'numeric', mu1 = 'integer', mu0 = 'integer', diff_exp = 'character')
r = read_regions(con = r_file, extraCols = extraCols, rename_score = 'coverage')

# Get premade CpG annotations
data('annotations', package = 'annotatr')

a = annotate_regions(

```

```
regions = r,
annotations = annotations,
ignore.strand = TRUE)
```

---

annotations	<i>example_annotations data</i>
-------------	---------------------------------

---

### Description

A GRanges of precomputed annotations for CpG features. Created by doing `build_annotations(genome='hg19', anno = 'hg19_cpgs')`.

### Usage

```
annotations
```

### Format

A GRanges object with the CpG feature annotations for hg19 and containing mcols:

**id** The internal ID for the annotation

**tx\_id** All NA, since these are not associated with tx\_ids

**gene\_id** All NA, since there are not associated Entrez IDs

**symbols** All NA, since there are not associated gene symbols

**type** A character indicating the type of annotation. Including: 'hg19\_cpg\_islands', 'hg19\_cpg\_shores', 'hg19\_cpg\_shelves', and 'hg19\_cpg\_inter'.

### Source

The AnnotationHub resource for hg19 CpG features.

---

annotatr	<i>annotatr: Annotation of Genomic Regions to Functional Annotations</i>
----------	--

---

### Description

Given a set of genomic sites/regions (e.g. ChIP-seq peaks, CpGs, differentially methylated CpGs or regions, SNPs, etc.) it is often of interest to investigate the intersecting functional annotations. Such annotations include those relating to gene models (promoters, 5'UTRs, exons, introns, and 3'UTRs), CpGs (CpG islands, CpG shores, CpG shelves), the non-coding genome, and enhancers. The `annotatr` package provides an easy way to summarize and visualize the intersection of genomic sites/regions with the above functional annotations.

---

annotatr_cache	<i>A global-variable to hold custom annotations loaded in an R session</i>
----------------	--

---

**Description**

Code thanks to Martin Morgan. This is a global variable that will store custom annotations that a user reads in during a session in which annotatr is loaded.

**Usage**

```
annotatr_cache
```

**Format**

An object of class list of length 3.

**Value**

An environment to contain custom annotations from read\_annotations.

**Examples**

```
# Example usage
annotatr_cache$set("foo", 1:10)
annotatr_cache$get("foo")

# Read in a BED3 file as a custom annotation
file = system.file('extdata', 'test_annotations_3.bed', package='annotatr')
# The custom annotation is added to the annotatr_cache environment in this function
read_annotations(con = file, name = 'test', genome = 'hg19')
# The result of read_annotations() is not visible in .GlobalEnv, instead
# need to use the get method
print(annotatr_cache$get('hg19_custom_test'))
# See what is in the annotatr_cache
annotatr_cache$list_env()
```

---

build_ah_annots	<i>A helper function to build arbitrary annotations from AnnotationHub</i>
-----------------	--

---

**Description**

A helper function to build arbitrary annotations from AnnotationHub

**Usage**

```
build_ah_annots(genome, ah_codes, annotation_class)
```

**Arguments**

genome	The genome assembly.
ah_codes	A named character vector giving the AnnotationHub accession number (e.g. AH23256), and whose name describes what the annotation is (e.g. Gm12878_H3K4me3).
annotation_class	A string to name the group of annotations in ah_codes

**Value**

A GRanges object stored in annotatr\_cache. To view an annotation built with this function, do `annotatr_cache$get(name)`. To add these annotations to a set of annotations, include `'[genome]_[annotation_class]'` in the call to `build_annotations()`. See example below.

**Examples**

```
# Create a named vector for the AnnotationHub accession codes with desired names
h3k4me3_code = c('Gm12878' = 'AH23256')
# Fetch ah_codes from AnnotationHub and create annotations annotatr understands
build_ah_annots(genome = 'hg19', ah_codes = h3k4me3_code, annotation_class = 'H3K4me3')
# The annotations as they appear in annotatr_cache
annot_name = c('hg19_H3K4me3_Gm12878')
# Build the annotations right before annotating any regions
annotations = build_annotations(genome = 'hg19', annotations = annot_name)
```

---

build_annotations	<i>A function to build annotations from TxDb.* and AnnotationHub resources</i>
-------------------	--

---

**Description**

Create a GRanges object consisting of all the desired annotations. Supported annotation codes are listed by `builtin_annotations()`. The basis for enhancer annotations are FANTOM5 data, the basis for CpG related annotations are CpG island tracks from AnnotationHub, and the basis for genic annotations are from the TxDb.\* and org.db group of packages.

**Usage**

```
build_annotations(genome, annotations)
```

**Arguments**

genome	The genome assembly.
annotations	A character vector of annotations to build. Valid annotation codes are listed with <code>builtin_annotations()</code> . The "basicgenes" shortcut builds the following regions: 1-5Kb upstream of TSSs, promoters, 5UTRs, exons, introns, and 3UTRs. The "cpgs" shortcut builds the following regions: CpG islands, shores, shelves, and interCGI regions. NOTE: Shortcuts need to be appended by the genome, e.g. hg19_basicgenes. Custom annotations whose names are of the form <code>[genome]_custom_[name]</code> should also be included. Custom annotations should be read in and converted to GRanges with <code>read_annotations()</code> . They can be for a supported_genome(), or for an unsupported genome.

**Value**

A GRanges object of all the annotations combined. The mcols are `id`, `tx_id`, `gene_id`, `symbol`, `type`. The `id` column is a unique name, the `tx_id` column is either a UCSC knownGene transcript ID (genic annotations) or a Ensembl transcript ID (lncRNA annotations), the `gene_id` is the Entrez ID, the `symbol` is the gene symbol from the `org.*.eg.db` mapping from the Entrez ID, and the `type` is of the form `[genome]_[type]_[name]`.

**Examples**

```
# Example with hg19 gene promoters
annots = c('hg19_genes_promoters')
annots_gr = build_annotations(genome = 'hg19', annotations = annots)

# See vignette for an example with custom annotation
```

---

build_cpg_annots	<i>A helper function to build CpG related annotations.</i>
------------------	--

---

**Description**

Using the AnnotationHub package, extract CpG island track for the appropriate genome and construct the shores, shelves, and interCGI annotations as desired.

**Usage**

```
build_cpg_annots(
  genome = annotatr::builtin_genomes(),
  annotations = annotatr::builtin_annotations()
)
```

**Arguments**

<code>genome</code>	The genome assembly.
<code>annotations</code>	A character vector with entries of the form <code>[genome]_cpg_{islands,shores,shelves,inter}</code> .

**Value**

A list of GRanges objects.

---

build\_enhancer\_annots *A helper function to build enhancer annotations for hg19 and mm10 from FANTOM5.*

---

### Description

A helper function to build enhancer annotations for hg19 and mm10 from FANTOM5.

### Usage

```
build_enhancer_annots(genome = c("hg19", "hg38", "mm9", "mm10"))
```

### Arguments

genome            The genome assembly.

### Value

A GRanges object.

---

build\_gene\_annots *A helper function to build genic annotations.*

---

### Description

Using the TxDb.\* group of packages, construct genic annotations consisting of any combination of 1-5kb upstream of a TSS, promoters (< 1kb from TSS), 5UTRs, CDS, exons, first exons, introns, intron/exon and exon/intron boundaries, 3UTRs, and intergenic.

### Usage

```
build_gene_annots(
  genome = annotatr::builtin_genomes(),
  annotations = annotatr::builtin_annotations()
)
```

### Arguments

genome            The genome assembly.

annotations      A character vector with entries of the form [genome]\_genes\_{1 to 5kb,promoters,5UTRs,cds,exon}

### Value

A list of GRanges objects with unique id of the form [type]:i, tx\_id being the UCSC knownGene transcript name, gene\_id being the Entrez Gene ID, symbol being the gene symbol from the Entrez ID to symbol mapping in org.db for that species, and type being the annotation type.



---

build_hmm_annots	<i>A helper function to build chromHMM annotations for hg19 from UCSC Genome Browser.</i>
------------------	---

---

**Description**

A helper function to build chromHMM annotations for hg19 from UCSC Genome Browser.

**Usage**

```
build_hmm_annots(  
  genome = c("hg19"),  
  annotations = annotatr::builtin_annotatins()  
)
```

**Arguments**

genome	The genome assembly.
annotations	A character vector of valid chromatin state annotatin codes.

**Value**

A GRanges object.

---

build_lncrna_annots	<i>A helper function to build lncRNA annotations.</i>
---------------------	---

---

**Description**

Using the AnnotationHub package, retrieve transcript level lncRNA annotations for either human (GRCh38) or mouse (GRCm38). If the genome is 'hg19', use the permalink from GENCODE and `rtracklayer::import()` to download and process.

**Usage**

```
build_lncrna_annots(genome = c("hg19", "hg38", "mm10"))
```

**Arguments**

genome	The genome assembly.
--------	----------------------

**Value**

A GRanges object with `id` giving the transcript\_type from the GENCODE file, `tx_id` being the Ensembl transcript ID, `gene_id` being the Entrez ID coming from a mapping of gene symbol to Entrez ID, `symbol` being the gene\_name from the GENCODE file, and the type being `[genome]_lncrna_gencode`.

---

builtin\_annotatons     *Function listing which annotations are available.*

---

**Description**

This includes the shortcuts. The `expand_annotatons()` function helps handle the shortcuts.

**Usage**

```
builtin_annotatons()
```

**Value**

A character vector of available annotations.

**Examples**

```
builtin_annotatons()
```

---

builtin\_genomes     *Function returning supported TxDb.\* genomes*

---

**Description**

Function returning supported TxDb.\* genomes

**Usage**

```
builtin_genomes()
```

**Value**

A character vector of genomes for supported TxDb.\* packages

**Examples**

```
builtin_genomes()
```

---

check_annotations	<i>Function to check for valid annotations</i>
-------------------	--

---

**Description**

Gives errors if any annotations are not in builtin\_annotations() (and they are not in the required custom format), basicgenes are used, or the genome prefixes are not the same for all annotations.

**Usage**

```
check_annotations(annotations)
```

**Arguments**

annotations      A character vector of annotations possibly using the shortcuts

**Value**

If all the checks on the annotations pass, returns NULL to allow code to move forward.

---

expand_annotations	<i>Function to expand annotation shortcuts</i>
--------------------	--

---

**Description**

Function to expand annotation shortcuts

**Usage**

```
expand_annotations(annotations)
```

**Arguments**

annotations      A character vector of annotations, possibly using the shortcut accessors

**Value**

A vector of data accession-ized names that are ordered from upstream to downstream in the case of knownGenes and islands to interCGI in the case of cpGs.

get\_cellline\_from\_code

*Function to return cell line from chromatin annotation code*

---

### **Description**

Function to return cell line from chromatin annotation code

### **Usage**

```
get_cellline_from_code(code)
```

### **Arguments**

code            The annotation code, used in build\_annotations().

### **Value**

A string of the cell line used in a chromatin annotation code

---

get\_cellline\_from\_shortcut

*Function to return cell line from chromatin annotation shortcut*

---

### **Description**

Function to return cell line from chromatin annotation shortcut

### **Usage**

```
get_cellline_from_shortcut(shortcut)
```

### **Arguments**

shortcut        The annotation shortcut, used in build\_annotations().

### **Value**

A string of the cell line used in a chromatin annotation shortcut

---

get_orgdb_name	<i>Function to get correct org.* package name based on genome</i>
----------------	---

---

**Description**

Function to get correct org.\* package name based on genome

**Usage**

```
get_orgdb_name(genome = annotatr::builtin_genomes())
```

**Arguments**

genome            A string giving the genome assembly.

**Value**

A string giving the correct org for org.db packages. e.g. hg19 -> Hs.

---

get_txdb_name	<i>Function to get correct TxDb.* package name based on genome</i>
---------------	--

---

**Description**

Function to get correct TxDb.\* package name based on genome

**Usage**

```
get_txdb_name(genome = annotatr::builtin_genomes())
```

**Arguments**

genome            A string giving the genome assembly.

**Value**

A string giving the name of the correct TxDb.\* package name based on genome.

---

plot\_annotation      *Plot the number of regions per annotation*

---

### Description

Given a GRanges of annotated regions, plot the number of regions with the corresponding genomic annotations used in annotation\_order. If a region is annotated to multiple annotations of the same annot. type, the region will only be counted once in the corresponding bar plot. For example, if a region were annotated to multiple exons, it would only count once toward the exon bar in the plot, but if it were annotated to an exon and an intron, it would count towards both.

### Usage

```
plot_annotation(
  annotated_regions,
  annotated_random,
  annotation_order = NULL,
  plot_title,
  x_label,
  y_label,
  quiet = FALSE
)
```

### Arguments

annotated\_regions      The GRanges result of annotate\_regions().

annotated\_random      The GRanges result of annotate\_regions() on the randomized regions created from randomize\_regions().

annotation\_order      A character vector which doubles as the subset of annotations desired for the plot as well as the ordering. If NULL, all annotations are displayed.

plot\_title      A string used for the title of the plot. If missing, no title is displayed.

x\_label      A string used for the x-axis label. If missing, no x-axis label is displayed.

y\_label      A string used for the y-axis label. If missing, no y-axis label is displayed.

quiet      Print progress messages (FALSE) or not (TRUE).

### Value

A ggplot object which can be viewed by calling it, saved with ggplot2::ggsave, or edited.

### Examples

```
#####
# An example of ChIP-seq peaks with signalValue used for score

# Get premade CpG annotations
data('annotations', package = 'annotatr')
```

```
chip_bed = system.file('extdata', 'Gm12878_Stat3_chr2.bed.gz', package = 'annotatr')
chip_regions = read_regions(con = chip_bed, genome = 'hg19')

chip_rnd = randomize_regions(regions = chip_regions)

chip_annots = annotate_regions(
  regions = chip_regions,
  annotations = annotations,
  ignore.strand = TRUE)

chip_rnd_annots = annotate_regions(
  regions = chip_rnd,
  annotations = annotations,
  ignore.strand = TRUE)

annots_order = c(
  'hg19_cpg_islands',
  'hg19_cpg_shores')

p_annots = plot_annotation(annotated_regions = chip_annots,
  annotation_order = annots_order)
p_annots_rnd = plot_annotation(annotated_regions = chip_annots,
  annotated_random = chip_rnd_annots, annotation_order = annots_order)
```

---

plot\_categorical

*Plot a categorical data variable over another*

---

## Description

Given a GRanges of annotated regions from `annotate_regions()`, visualize the the distribution of categorical data fill in categorical data x. A bar representing the distribution of all fill in x will be added according to the contents of fill. This is the distribution over all values of x. Additionally, when `annotated_random` is not missing, a "Random Regions" bar shows the distribution of random regions over fill.

## Usage

```
plot_categorical(
  annotated_regions,
  annotated_random,
  x,
  fill = NULL,
  x_order = NULL,
  fill_order = NULL,
  position = "stack",
  plot_title,
  legend_title,
  x_label,
  y_label,
  quiet = FALSE
)
```

**Arguments**

annotated_regions	The GRanges result of <code>annotate_regions()</code> .
annotated_random	The GRanges result of <code>annotate_regions()</code> on the randomized regions created from <code>randomize_regions()</code> . Random regions can only be used with <code>fill == 'annot.type'</code> .
x	One of <code>'annot.type'</code> or a categorical data column, indicating whether annotation classes or data classes will appear on the x-axis.
fill	One of <code>'annot.type'</code> , a categorical data column, or NULL, indicating whether annotation classes or data classes will fill the bars. If NULL then the bars will be the total counts of the x classes.
x_order	A character vector that subsets and orders the x classes. Default NULL, uses existing values.
fill_order	A character vector that subsets and orders the fill classes. Default NULL, uses existing values.
position	A string which has the same possible values as in <code>ggplot2::geom_bar(..., position)</code> , i.e., <code>'stack'</code> , <code>'fill'</code> , <code>'dodge'</code> , etc.
plot_title	A string used for the title of the plot. If missing, no title is displayed.
legend_title	A string used for the legend title to describe fills (if fill is not NULL). Default displays corresponding variable name.
x_label	A string used for the x-axis label. If missing, corresponding variable name used.
y_label	A string used for the y-axis label. If missing, corresponding variable name used.
quiet	Print progress messages (FALSE) or not (TRUE).

**Details**

For example, if a differentially methylated region has the categorical label `hyper`, and is annotated to a promoter, a 5UTR, two exons, and an intron. Each annotation will appear in the All bar once. Likewise for the hyper bar if the differential methylation status is chosen as `x` with `annot.type` chosen as `fill`.

**Value**

A `ggplot` object which can be viewed by calling it, or saved with `ggplot2::ggsave`.

**Examples**

```
# Get premade CpG annotations
data('annotations', package = 'annotatr')

dm_file = system.file('extdata', 'IDH2mut_v_NBM_multi_data_chr9.txt.gz', package = 'annotatr')
extraCols = c(diff_meth = 'numeric', mu1 = 'numeric', mu0 = 'numeric')
dm_regions = read_regions(con = dm_file, extraCols = extraCols, genome = 'hg19',
  rename_score = 'pval', rename_name = 'DM_status', format = 'bed')
dm_regions = dm_regions[1:1000]

dm_annots = annotate_regions(
  regions = dm_regions,
  annotations = annotations,
```



```

        ignore.strand = TRUE)

dm_order = c(
  'hyper',
  'hypo')
cpg_order = c(
  'hg19_cpg_islands',
  'hg19_cpg_shores',
  'hg19_cpg_shelves',
  'hg19_cpg_inter')

dm_vn = plot_categorical(
  annotated_regions = dm_annots,
  x = 'DM_status',
  fill = 'annot.type',
  x_order = dm_order,
  fill_order = cpg_order,
  position = 'fill',
  legend_title = 'knownGene Annotations',
  x_label = 'DM status',
  y_label = 'Proportion')

# Create randomized regions
dm_rnd_regions = randomize_regions(regions = dm_regions)
dm_rnd_annots = annotate_regions(
  regions = dm_rnd_regions,
  annotations = annotations,
  ignore.strand = TRUE)

dm_vn_rnd = plot_categorical(
  annotated_regions = dm_annots,
  annotated_random = dm_rnd_annots,
  x = 'DM_status',
  fill = 'annot.type',
  x_order = dm_order,
  fill_order = cpg_order,
  position = 'fill',
  legend_title = 'knownGene Annotations',
  x_label = 'DM status',
  y_label = 'Proportion')

```

---

plot\_coannotations      *Plot pair-wise annotations across regions*

---

### Description

All co-occurring annotations associated with a region are computed and displayed as a heatmap.

### Usage

```

plot_coannotations(
  annotated_regions,
  annotation_order = NULL,

```

```

    plot_title,
    axes_label,
    quiet = FALSE
  )

```

### Arguments

annotated_regions	The GRanges result of <code>annotate_regions()</code> .
annotation_order	A character vector which doubles as the subset of annotations desired for plot as well as the ordering. If NULL, all annotations are displayed.
plot_title	A string used for the title of the plot. If missing, no plot title label is displayed.
axes_label	A string used for the axis labels. If missing, corresponding variable name used.
quiet	Print progress messages (FALSE) or not (TRUE).

### Details

As with `plot_annotation()`, the number in each cell is the number of unique regions annotated to the pair of annotations.

For example, if a region is annotated to both a CpG shore and to two different exons simultaneously, the region will only be counted once in the CpG shore / exon cell. NOTE, this same region will count once in both the CpG shore and exon cells on the diagonal.

### Value

A ggplot object which can be viewed by calling it, saved with `ggplot2::ggsave`, or edited.

### Examples

```

# Get premade CpG annotations
data('annotations', package = 'annotatr')

dm_file = system.file('extdata', 'IDH2mut_v_NBM_multi_data_chr9.txt.gz', package = 'annotatr')
extraCols = c(diff_meth = 'numeric', mu1 = 'numeric', mu0 = 'numeric')
dm_regions = read_regions(con = dm_file, extraCols = extraCols,
  rename_score = 'pval', rename_name = 'DM_status', format = 'bed')
dm_regions = dm_regions[1:1000]

dm_annots = annotate_regions(
  regions = dm_regions,
  annotations = annotations,
  ignore.strand = TRUE)

all_order = c(
  'hg19_cpg_islands',
  'hg19_cpg_shores',
  'hg19_cpg_shelves',
  'hg19_cpg_inter')

dm_vs_ca = plot_coannotations(
  annotated_regions = dm_annots,
  annotation_order = all_order,
  axes_label = 'Annotations',

```

```
plot_title = 'Co-occurrence of Annotations')
```

---

plot_numerical	<i>Plot numerical data over regions or regions summarized over annotations</i>
----------------	--

---

## Description

This function produces either histograms over facet, or x-y scatterplots over facet. In the case of histograms over facets, the All distribution (hollow histogram with red outline) is the distribution of x over all the regions in the data. The facet specific distributions (solid gray) are the distribution of x over the regions in each facet. For example, a CpG with associated percent methylation annotated to a CpG island and a promoter will count once in the All distribution, but will count once each in the CpG island and promoter facet distributions.

## Usage

```
plot_numerical(
  annotated_regions,
  x,
  y,
  facet,
  facet_order,
  bin_width = 10,
  plot_title,
  x_label,
  y_label,
  legend_facet_label,
  legend_cum_label,
  quiet = FALSE
)
```

## Arguments

annotated_regions	A GRanges returned from <code>annotate_regions()</code> . If the data is not summarized, the data is at the region level. If it is summarized, it represents the average or standard deviation of the regions by the character vector used for by in <code>summarize_numerical()</code> .
x	A string indicating the column of the GRanges to use for the x-axis.
y	A string indicating the column of the GRanges to use for the y-axis. If missing, a histogram over x will be plotted. If not missing, a scatterplot is plotted.
facet	A string, or character vector of two strings, indicating indicating which categorical variable(s) in the GRanges to make <code>ggplot2</code> facets over. When two facets are given, the first entry is the vertical facet and the second entry is the horizontal facet. Default is <code>annot.type</code> .
facet_order	A character vector, or list of character vectors if <code>facet</code> has length 2, which gives the order of the facets, and can be used to subset the column in the GRanges used for the facet. For example, if <code>facet = 'annot.type'</code> , then the annotations maybe subsetted to just CpG annotations. Default is <code>NULL</code> , meaning all annotations in their default order are used.

bin_width	An integer indicating the bin width of the histogram used for score. Default 10. Select something appropriate for the data. NOTE: This is only used if y is NULL.
plot_title	A string used for the title of the plot. If missing, no title is displayed.
x_label	A string used for the x-axis label. If missing, no x-axis label is displayed.
y_label	A string used for the y-axis label. If missing, no y-axis label is displayed.
legend_facet_label	A string used to label the gray bar portion of the legend. Defaults to "x in facet".
legend_cum_label	A string used to label the red outline portion of the legend. Defaults to "All in x".
quiet	Print progress messages (FALSE) or not (TRUE).

### Value

A ggplot object which can be viewed by calling it, or saved with `ggplot2::ggsave`.

### Examples

```
# An example with multi-columned data

# Get premade CpG annotations
data('annotations', package = 'annotatr')

dm_file = system.file('extdata', 'IDH2mut_v_NBM_multi_data_chr9.txt.gz', package = 'annotatr')
extraCols = c(diff_meth = 'numeric', mu1 = 'numeric', mu0 = 'numeric')
dm_regions = read_regions(con = dm_file, extraCols = extraCols,
  rename_score = 'pval', rename_name = 'DM_status', format = 'bed')
dm_regions = dm_regions[1:1000]

# Annotate the regions
dm_annots = annotate_regions(
  regions = dm_regions,
  annotations = annotations,
  ignore.strand = TRUE)

# Plot histogram of group 1 methylation rates across the CpG annotations.
# NOTE: Overall distribution (everything in \code{facet_order})
# is plotted in each facet for comparison.
dm_vs_regions_mu1 = plot_numerical(
  annotated_regions = dm_annots,
  x = 'mu1',
  facet = 'annot.type',
  facet_order = c('hg19_cpg_islands', 'hg19_cpg_shores',
    'hg19_cpg_shelves', 'hg19_cpg_inter'),
  bin_width = 5,
  plot_title = 'Group 1 Methylation over CpG Annotations',
  x_label = 'Group 1 Methylation')

# Plot histogram of group 1 methylation rates across the CpG annotations
# crossed with DM_status
dm_vs_regions_diffmeth = plot_numerical(
  annotated_regions = dm_annots,
  x = 'diff_meth',
  facet = c('annot.type', 'DM_status'),
```

```

facet_order = list(c('hg19_genes_promoters', 'hg19_genes_5UTRs', 'hg19_cpg_islands'), c('hyper', 'hypo', 'n
  bin_width = 5,
  plot_title = 'Group 0 Region Methylation In Genes',
  x_label = 'Methylation Difference')

# Can also use the result of annotate_regions() to plot two numerical
# data columns against each other for each region, and facet by annotations.
dm_vs_regions_annot = plot_numerical(
  annotated_regions = dm_annots,
  x = 'mu0',
  y = 'mu1',
  facet = 'annot.type',
  facet_order = c('hg19_cpg_islands', 'hg19_cpg_shores',
    'hg19_cpg_shelves', 'hg19_cpg_inter'),
  plot_title = 'Region Methylation: Group 0 vs Group 1',
  x_label = 'Group 0',
  y_label = 'Group 1')

# Another example, but using differential methylation status as the facets.
dm_vs_regions_name = plot_numerical(
  annotated_regions = dm_annots,
  x = 'mu0',
  y = 'mu1',
  facet = 'DM_status',
  facet_order = c('hyper', 'hypo', 'none'),
  plot_title = 'Region Methylation: Group 0 vs Group 1',
  x_label = 'Group 0',
  y_label = 'Group 1')

```

---

plot\_numerical\_coannotations

*Plot numerical data occurring in pairs of annotations*

---

## Description

Plot numerical data associated with regions occurring in annot1, annot2 and in both. As with plot\_numerical(), the result is a plot of histograms or x-y scatterplots.

## Usage

```

plot_numerical_coannotations(
  annotated_regions,
  x,
  y,
  annot1,
  annot2,
  bin_width = 10,
  plot_title,
  x_label,
  y_label,
  legend_facet_label,
  legend_cum_label,

```

```

    quiet = FALSE
  )

```

### Arguments

annotated_regions	A GRanges returned from <code>annotate_regions()</code> .
x	A string indicating the column of the GRanges to use for the x-axis.
y	A string indicating the column of the GRanges to use for the y-axis. If missing, a histogram over x will be plotted. If not missing, a scatterplot is plotted.
annot1	A string indicating the first annotation type.
annot2	A string indicating the second annotation type.
bin_width	An integer indicating the bin width of the histogram used for score. Default 10. Select something appropriate for the data. NOTE: This is only used if y is NULL.
plot_title	A string used for the title of the plot. If missing, no title is displayed.
x_label	A string used for the x-axis label. If missing, no x-axis label is displayed.
y_label	A string used for the y-axis label. If missing, no y-axis label is displayed.
legend_facet_label	A string used to label the gray bar portion of the legend. Defaults to "x in annot pair".
legend_cum_label	A string used to label the red outline portion of the legend. Defaults to "All x".
quiet	Print progress messages (FALSE) or not (TRUE).

### Details

For example, a CpG with associated percent methylation annotated to a CpG island and a promoter will count once in the All distribution and once in the CpG island / promoter facet distribution. However, a CpG associated only with a promoter will count once in the All distribution and once in the promoter / promoter distribution.

### Value

A ggplot object which can be viewed by calling it, or saved with `ggplot2::ggsave`.

### Examples

```

# Get premade CpG annotations
data('annotations', package = 'annotatr')

dm_file = system.file('extdata', 'IDH2mut_v_NBM_multi_data_chr9.txt.gz', package = 'annotatr')
extraCols = c(diff_meth = 'numeric', mu1 = 'numeric', mu0 = 'numeric')
dm_regions = read_regions(con = dm_file, extraCols = extraCols,
  rename_score = 'pval', rename_name = 'DM_status', format = 'bed')
dm_regions = dm_regions[1:1000]

dm_annots = annotate_regions(
  regions = dm_regions,
  annotations = annotations,
  ignore.strand = TRUE)

```

```
dm_vs_num_co = plot_numerical_coannotations(
  annotated_regions = dm_annots,
  x = 'mu0',
  annot1 = 'hg19_cpg_islands',
  annot2 = 'hg19_cpg_shelves',
  bin_width = 5,
  plot_title = 'Group 0 Perc. Meth. in CpG Islands and Promoters',
  x_label = 'Percent Methylation')
```

---

randomize_regions	<i>Randomize Regions</i>
-------------------	--------------------------

---

## Description

randomize\_regions is a wrapper function for `regioneR::randomizeRegions()` that simplifies the creation of randomized regions for an input set of regions read with `read_regions()`. It relies on the `seqlengths` of regions in order to build the appropriate genome object for `regioneR::randomizeRegions()`.

## Usage

```
randomize_regions(
  regions,
  allow.overlaps = TRUE,
  per.chromosome = TRUE,
  quiet = FALSE
)
```

## Arguments

regions	A GRanges object from <code>read_regions</code> .
allow.overlaps	A logical stating whether random regions can overlap input regions (TRUE) or not (FALSE). Default TRUE.
per.chromosome	A logical stating whether the random regions should remain on the same chromosome (TRUE) or not (FALSE). Default TRUE.
quiet	Print progress messages (FALSE) or not (TRUE).

## Details

NOTE: The data associated with the input regions are not passed on to the random regions.

## Value

A GRanges object of randomized regions based on regions from `read_regions()`. NOTE: Data associated with the original regions is not attached to the randomized regions.

## Examples

```
# Create random region set based on ENCODE ChIP-seq data
file = system.file('extdata', 'Gm12878_Stat3_chr2.bed.gz', package = 'annotatr')
r = read_regions(con = file, genome = 'hg19')

random_r = randomize_regions(regions = r)
```

---

read_annotations	<i>Read custom annotations</i>
------------------	--------------------------------

---

### Description

read\_annotations() is a wrapper for the rtracklayer::import() function that creates a GRanges object matching the structure of annotations built with build\_annotations(). The structure is defined by GRanges, with the mcols() with names c('id', 'gene\_id', 'symbol', 'type').

### Usage

```
read_annotations(con, name, genome = NA, format, extraCols = character(), ...)
```

### Arguments

con	A path, URL, connection or BEDFile object. See rtracklayer::import.bed() documentation.
name	A string for the name of the annotations to be used in the name of the object, [genome]_custom_[name]
genome	From rtracklayer::import(): The identifier of a genome, or NA if unknown. Typically, this is a UCSC identifier like 'hg19'. An attempt will be made to derive the seqinfo on the return value using either an installed BSgenome package or UCSC, if network access is available.
format	From rtracklayer::import(): The format of the output. If not missing, should be one of 'bed', 'bed15', 'bedGraph' or 'bedpe'. If missing and 'con' is a filename, the format is derived from the file extension. This argument is unnecessary when 'con' is a derivative of 'RTLFile'.
extraCols	From rtracklayer::import.bed(): A character vector in the same form as 'colClasses' from 'read.table'. It should indicate the name and class of each extra/special column to read from the BED file. As BED does not encode column names, these are assumed to be the last columns in the file. This enables parsing of the various BEDX+Y formats.
...	Parameters to pass onto the format-specific method of rtracklayer::import().

### Value

A GRanges object stored in annotatr\_cache. To view a custom annotation, do annotatr\_cache\$get(name). To add a custom annotation to the set of annotations, include '[genome]\_custom\_[name]' in the call to build\_annotations(). See example below.

### Examples

```
# Read in a BED3 file as a custom annotation
file = system.file('extdata', 'test_annotations_3.bed', package='annotatr')
read_annotations(con = file, name = 'test', genome = 'hg19')
build_annotations(genome = 'hg19', annotations = 'hg19_custom_test')

print(annotatr_cache$get('hg19_custom_test'))
```



---

read_regions	<i>Read genomic regions in BEDX+Y format</i>
--------------	--

---

### Description

read\_regions() reads genomic regions by calling the rtracklayer::import() function. This function can automatically deal with BEDX files from BED3 to BED6. For BED6+Y, the extraCols argument should be used to correctly interpret the extra columns.

### Usage

```
read_regions(
  con,
  genome = NA,
  format,
  extraCols = character(),
  rename_name,
  rename_score,
  ...
)
```

### Arguments

con	A path, URL, connection or BEDFile object. See rtracklayer::import() documentation.
genome	From rtracklayer::import(): The identifier of a genome, or NA if unknown. Typically, this is a UCSC identifier like 'hg19'. An attempt will be made to derive the seqinfo on the return value using either an installed BSgenome package or UCSC, if network access is available.
format	From rtracklayer::import(): The format of the output. If not missing, should be one of 'bed', 'bed15', 'bedGraph' or 'bedpe'. If missing and 'con' is a filename, the format is derived from the file extension. This argument is unnecessary when 'con' is a derivative of 'RTLFile'.
extraCols	From rtracklayer::import(): A character vector in the same form as 'colClasses' from 'read.table'. It should indicate the name and class of each extra/special column to read from the BED file. As BED does not encode column names, these are assumed to be the last columns in the file. This enables parsing of the various BEDX+Y formats.
rename_name	A string to rename the name column of the BED file. For example, if the name column actually contains a categorical variable.
rename_score	A string to rename the score column of the BED file. For example, if the score column represents a quantity about the data besides the score in the BED specification.
...	Parameters to pass onto the format-specific method of rtracklayer::import().

### Details

NOTE: The name (4th) and score (5th) columns are so named. If these columns have a particular meaning for your data, they should be renamed with the rename\_name and/or rename\_score parameters.

**Value**

A GRanges object.

**Examples**

```
# Example of reading a BED6+3 file where the last 3 columns are non-standard
file = system.file('extdata', 'IDH2mut_v_NBM_multi_data_chr9.txt.gz', package = 'annotatr')
extraCols = c(diff_meth = 'numeric', mu0 = 'numeric', mu1 = 'numeric')
gr = read_regions(con = file, genome = 'hg19', extraCols = extraCols, format = 'bed',
  rename_name = 'DM_status', rename_score = 'pval')
```

---

reformat\_hmm\_codes      *Function to recode classes from chromHMM type column*

---

**Description**

Function to recode classes from chromHMM type column

**Usage**

```
reformat_hmm_codes(hmm_codes)
```

**Arguments**

hmm\_codes      in the original form from UCSC Genome Browser track.

**Value**

A character vector of chromHMM classes with numbers and underscores removed.

---

subset\_order\_tbl      *Function to subset a tbl\_df or grouped\_df by a column*

---

**Description**

Function to subset a tbl\_df or grouped\_df by a column

**Usage**

```
subset_order_tbl(tbl, col, col_order)
```

**Arguments**

tbl      A tbl\_df or grouped\_df.  
col      A string indicating which column of of tbl to subset and order  
col\_order      A character vector indicating the order of col.

**Value**

A modified version of summary with col subsetted by col\_order.

---

summarize\_annotatations *Summarize annotation counts*

---

### Description

Given a GRanges of annotated regions, count the number of regions in each annotation type. If annotated\_random is not NULL, then the same is computed for the random regions.

### Usage

```
summarize_annotatations(annotated_regions, annotated_random, quiet = FALSE)
```

### Arguments

annotated\_regions  
The GRanges result of `annotate_regions()`.

annotated\_random  
The GRanges result of `annotate_regions()` on the randomized regions created from `randomize_regions()`.

quiet  
Print progress messages (FALSE) or not (TRUE).

### Details

If a region is annotated to multiple annotations of the same `annot.type`, the region will only be counted once. For example, if a region were annotated to multiple exons, it would only count once toward the exons, but if it were annotated to an exon and an intron, it would count towards both.

### Value

A `tbl_df` of the number of regions per annotation type.

### Examples

```
### An example of ChIP-seq peaks with signalValue

# Get premade CpG annotations
data('annotations', package = 'annotatr')

file = system.file('extdata', 'Gm12878_Stat3_chr2.bed.gz', package = 'annotatr')
r = read_regions(con = file, genome = 'hg19')

a = annotate_regions(
  regions = r,
  annotations = annotations,
  ignore.strand = TRUE,
  quiet = FALSE)

rnd = randomize_regions(regions = r)

rnd_annots = annotate_regions(
  regions = rnd,
  annotations = annotations,
  ignore.strand = TRUE,
```

```

quiet = FALSE)

# Summarize the annotated regions without randomized regions
s = summarize_annotations(annotated_regions = a)

# Summarize the annotated regions with randomized regions
s_rnd = summarize_annotations(
  annotated_regions = a,
  annotated_random = rnd_annots)

```

---

summarize\_categorical *Summarize categorical data over groupings of annotated regions*

---

## Description

Given a GRanges of annotated regions, count the number of regions when the annotations are grouped by categorical columns.

## Usage

```

summarize_categorical(
  annotated_regions,
  by = c("annot.type", "annot.id"),
  quiet = FALSE
)

```

## Arguments

annotated_regions	The GRanges result of <code>annotate_regions()</code> .
by	A character vector to group the data in <code>as.data.frame(annotated_regions)</code> by and tally over. Default is <code>c('annot.type', 'annot.id')</code> .
quiet	Print progress messages (FALSE) or not (TRUE).

## Details

If a region is annotated to multiple annotations of the same `annot.type`, the region will only be counted once. For example, if a region were annotated to multiple exons, it would only count once toward the exons, but if it were annotated to an exon and an intron, it would count towards both.

## Value

A grouped `dplyr::tbl_df` of the counts of groupings according to the `by` vector.

## Examples

```

# Get premade CpG annotations
data('annotations', package = 'annotatr')

r_file = system.file('extdata', 'test_read_multiple_data_nohead.bed', package='annotatr')
extraCols = c(pval = 'numeric', mu1 = 'integer', mu0 = 'integer', diff_exp = 'character')

```

```

r = read_regions(con = r_file, genome = 'hg19', extraCols = extraCols, rename_score = 'coverage')

a = annotate_regions(
  regions = r,
  annotations = annotations,
  ignore.strand = TRUE)

sc = summarize_categorical(
  annotated_regions = a,
  by = c('annot.type', 'name'),
  quiet = FALSE)

```

---

summarize\_numerical     *Summarize numerical data over groupings of annotated regions*

---

## Description

Given a GRanges of annotated regions, summarize numerical data columns based on a grouping.

## Usage

```

summarize_numerical(
  annotated_regions,
  by = c("annot.type", "annot.id"),
  over,
  quiet = FALSE
)

```

## Arguments

annotated_regions	The GRanges result of <code>annotate_regions()</code> .
by	A character vector of the columns of <code>as.data.frame(annotated_regions)</code> to group over. Default is <code>c(annot.type, annot.id)</code> .
over	A character vector of the numerical columns in <code>as.data.frame(annotated_regions)</code> to count, take the mean, and take the sd over after grouping according to the by column. NOTE: If more than one value is used, the naming scheme for the resulting <code>dplyr::tbl</code> summary columns are <code>COLNAME_n</code> , <code>COLNAME_mean</code> , <code>COLNAME_sd</code> . If over has length one, then the column names are n, mean, sd.
quiet	Print progress messages (FALSE) or not (TRUE).

## Details

NOTE: We do not take the distinct values of `seqnames`, `start`, `end`, `annot.type` as in the other `summarize_*()` functions because in the case of a region that intersected two distinct exons, using `distinct_()` would destroy the information of the mean of the numerical column over one of the exons, which is not desirable.

## Value

A grouped `dplyr::tbl_df`, and the count, mean, and sd of the cols by the groupings.

**Examples**

```

### Test on a very simple bed file to demonstrate different options

# Get premade CpG annotations
data('annotations', package = 'annotatr')

r_file = system.file('extdata', 'test_read_multiple_data_nohead.bed', package='annotatr')
extraCols = c(pval = 'numeric', mu1 = 'integer', mu0 = 'integer', diff_exp = 'character')
r = read_regions(con = r_file, genome = 'hg19', extraCols = extraCols, rename_score = 'coverage')

a = annotate_regions(
  regions = r,
  annotations = annotations,
  ignore.strand = TRUE)

# Testing over normal by
sn1 = summarize_numerical(
  annotated_regions = a,
  by = c('annot.type', 'annot.id'),
  over = c('coverage', 'mu1', 'mu0'),
  quiet = FALSE)

# Testing over a different by
sn2 = summarize_numerical(
  annotated_regions = a,
  by = c('diff_exp'),
  over = c('coverage', 'mu1', 'mu0'))

```

---

tidy_annotations	<i>Function to tidy up annotation accessors for visualization</i>
------------------	---

---

**Description**

Function to tidy up annotation accessors for visualization

**Usage**

```
tidy_annotations(annotations)
```

**Arguments**

`annotations` A character vector of annotations, in the order they are to appear in the visualization.

**Value**

A list of mappings from original annotation names to names ready for visualization.

# Index

## \* datasets

- annotations, [4](#)
- annotatr\_cache, [5](#)

- annotate\_regions, [3](#)
- annotations, [4](#)
- annotatr, [4](#)
- annotatr\_cache, [5](#)

- build\_ah\_annots, [5](#)
- build\_annotations, [6](#)
- build\_cpg\_annots, [7](#)
- build\_enhancer\_annots, [8](#)
- build\_gene\_annots, [8](#)
- build\_hmm\_annots, [9](#)
- build\_lncrna\_annots, [9](#)
- builtin\_annotations, [10](#)
- builtin\_genomes, [10](#)

- check\_annotations, [11](#)

- expand\_annotations, [11](#)

- get\_cellline\_from\_code, [12](#)
- get\_cellline\_from\_shortcut, [12](#)
- get\_orgdb\_name, [13](#)
- get\_txdb\_name, [13](#)

- plot\_annotation, [14](#)
- plot\_categorical, [15](#)
- plot\_coannotations, [17](#)
- plot\_numerical, [19](#)
- plot\_numerical\_coannotations, [21](#)

- randomize\_regions, [23](#)
- read\_annotations, [24](#)
- read\_regions, [25](#)
- reformat\_hmm\_codes, [26](#)

- subset\_order\_tbl, [26](#)
- summarize\_annotations, [27](#)
- summarize\_categorical, [28](#)
- summarize\_numerical, [29](#)

- tidy\_annotations, [30](#)