

Package ‘Platypus’

October 19, 2021

Type Package

Title Single-Cell Immune Repertoire and Gene Expression Analysis

Description We present 'Platypus', an open-source software platform providing a user-friendly interface to investigate B-cell receptor and T-cell receptor repertoires from scSeq experiments. 'Platypus' provides a framework to automate and ease the analysis of single-cell immune repertoires while also incorporating transcriptional information involving unsupervised clustering, gene expression and gene ontology (Yermanos A, et al (2021) <[doi:10.1093/nargab/lqab023](https://doi.org/10.1093/nargab/lqab023)>).

Version 3.2.3

Date 2021-09-23

Maintainer Alexander Yermanos <ayermanos@gmail.com>

Depends R(>= 3.5.0)

Imports BiocGenerics, Biostrings, dplyr, ggplot2, knitr, jsonlite, Matrix, plyr, reshape2, rmarkdown, seqinr, Seurat, SeuratObject, stringr, tibble, tidyr, utils, useful

Suggests AnnotationDbi, ape, circlize, cowplot, data.table, doParallel, fda.usc, fgsea, ggrepel, ggridges, ggseqlogo, grid, gridExtra, harmony, igraph, IRanges, limma, keras, methods, msgidbr, org.Mm.eg.db, pheatmap, phytools, purrr, scales, SummarizedExperiment, stringdist, tensorflow, vegan, testthat (>= 3.0.0)

License GPL-2

Encoding UTF-8

RoxygenNote 7.1.1

VignetteBuilder knitr

LazyData true

Config/testthat/edition 3

NeedsCompilation no

Author Alexander Yermanos [aut, cre],
Andreas Agrafiotis [ctb],
Raphael Kuhn [ctb],

Victor Kreiner [ctb],
 Damiano Robbiani [ctb],
 Josephine Yates [ctb],
 Jiami Han [ctb],
 Chrysa Papadopoulou [ctb],
 Cédric Weber [ctb],
 Danielle Shlesinger [ctb],
 Raphael Dizerens [ctb]

Repository CRAN

Date/Publication 2021-10-19 07:00:02 UTC

R topics documented:

AbForests_AntibodyForest	4
AbForests_CompareForests	7
AbForests_ConvertStructure	10
AbForests_CsvToDf	11
AbForests_ForestMetrics	11
AbForests_PlotGraphs	14
AbForests_PlyloToMatrix	15
AbForests_RemoveNets	16
AbForests_SubRepertoiresByCells	18
AbForests_SubRepertoiresByUniqueSeq	19
AbForests_UniqueAntibodyVariants	21
automate_GEX	22
Bcell_sequences_example_tree	24
Bcell_tree_2	25
call_MIXCR	25
class_switch_prob_hum	26
class_switch_prob_mus	27
clonofreq	27
clonofreq.isotype.data	28
clonofreq.isotype.plot	28
clonofreq.trans.data	29
clonofreq.trans.plot	30
cluster.id.igraph	30
colors	31
Echidna_vae_generate	32
get.avr.mut.data	33
get.avr.mut.plot	33
get.barplot.errorbar	34
get.elbow	35
get.n.node.data	35
get.n.node.plot	36
get.seq.distance	36
get.umap	37
get.vgu.matrix	37

GEX_automate	38
GEX_clonotype	40
GEX_cluster_genes	41
GEX_cluster_genes_heatmap	42
GEX_cluster_membership	43
GEX_coexpression_coefficient	44
GEX_DEgenes	45
GEX_DEgenes_persample	48
GEX_dottile_plot	49
GEX_GOterm	50
GEX_GSEA	52
GEX_heatmap	54
GEX_pairwise_DEGs	55
GEX_phenotype	56
GEX_phenotype_per_clone	57
GEX_proportions_barplot	58
GEX_scatter_coexpression	59
GEX_topN_DE_genes_per_cluster	60
GEX_visualize_clones	61
GEX_volcano	62
hotspot_df	64
hum_b_h	65
hum_b_l	65
hum_t_h	66
hum_t_l	67
iso_SHM_prob	67
mus_b_h	68
mus_b_l	68
mus_b_trans	69
mus_t_h	70
mus_t_l	71
new	71
no.empty.node	72
one_spot_df	72
pheno_SHM_prob	73
PlatypusDB_AIRR_to_VGM	74
PlatypusDB_fetch	75
PlatypusDB_find_CDR3s	77
PlatypusDB_list_projects	78
PlatypusDB_load_from_disk	79
PlatypusDB_VGM_to_AIRR	80
select.top	82
simulate_repertoire	83
small_vgm	87
special_v	88
trans_switch_prob_b	88
trans_switch_prob_t	89
umap.top.highlight	89

VDJ_alpha_beta_Vgene_circos	90
VDJ_analyze	91
VDJ_assemble_for_PnP	93
VDJ_call_MIXCR	95
VDJ_circos	96
VDJ_clonal_donut	97
VDJ_clonal_expansion	99
VDJ_clonal_lineages	101
VDJ_clonotype	102
VDJ_clonotype_clusters_circos	104
VDJ_contigs_to_vgm	106
VDJ_diversity	106
VDJ_dublets	108
VDJ_extract_germline	109
VDJ_GEX_clonal_lineage_clusters	110
VDJ_GEX_expansion	111
VDJ_GEX_integrate	112
VDJ_GEX_matrix	113
VDJ_GEX_overlay_clones	119
VDJ_GEX_stats	122
VDJ_isotypes_per_clone	123
vdj_length_prob	124
VDJ_logoplot_vector	126
VDJ_network	126
VDJ_overlap_heatmap	128
VDJ_per_clone	129
VDJ_plot_SHM	130
VDJ_reclonotype_list_arrange	132
VDJ_tree	133
VDJ_variants_per_clone	134
VDJ_Vgene_usage	135
VDJ_Vgene_usage_barplot	136
VDJ_Vgene_usage_stacked_barplot	137
VDJ_VJ_usage_circos	138

Index**140**

 AbForests_AntibodyForest

Infer and draw B cell evolutionary networks

Description

AntibodyForest takes the output of either ConvertStructure or CsvToDf or SubRepertoires or RemoveNets and outputs B cell phylogenetic networks in tree format. There is also the possibility to give the full-length list of clonal lineages, which contains both isotype and transcriptional cluster information, only when no prior data transformation is desired. Each network represents a clonal

lineage, referring to the number of B cell receptor sequences originating from an independent V(D)J recombination event. Each vertex represents a unique recovered full-length variable heavy and light chain antibody sequence of a clonal family. Edges separating nodes are drawn given that clonal variants are similarly related according to their Levenshtein distance. Edge weights are extracted from the distance matrix apart from the special case of unmutated germline, in which the weights of outgoing edges from it are either set to 1 or to the difference between the corresponding distance from the matrix and the absolute value of the difference between the sequence lengths of germline and corresponding connected nodes. At tree building, starting from the reference ancestral germline, each node is connected to nodes that can be reached via the minimum distance based on the distance matrix calculation. Therefore, potential edges that go back to previous tree layers along with bidirectional circles are eliminated. Polytomies, displayed by B cell clones producing multiple distinct offsprings, are resolved in case of reaching nodes with equal minimum distance. Indeed, the algorithm removes edges either randomly from the recipient nodes, based on the node closest or farthest from the germline, considering the number of intermediate nodes or edge path length, or the highest/lowest counting of cells on the present node. Additional ties are settled by random edge selection. Consequently, parsimony holds, meaning that each daughter node has only one parent. Distinct tree topologies enable to visually investigate the trade-off between balance and evolution, and further quantify the amount of diversification of the subsequent detected clonal abundant clones during somatic hypermutation and class switching. The minimum decision-based criterion determines the amount of balance presented in the tree, while the maximum decision-based method the amount of evolution presented in the tree. Single color or color distribution on each node demonstrates the proportion of B cells with the specific isotype(s) or transcriptional cluster(s), while setting the size of vertices can be performed based on the number of unique sequences per clone, vertex betweenness and vertex closeness. Scaling of nodes by their relative clonal expansion assists in pinpointing identical antibody sequences across a multitude of B cells. Node labeling can depict clonal frequency.

Usage

```
AbForests_AntibodyForest(  
  full_list,  
  csv,  
  files,  
  distance_mat,  
  clonal_frequency,  
  scaleByClonalFreq,  
  weight,  
  tie_flag,  
  scaleBybetweenness,  
  scaleBycloccloseness_metr,  
  opt,  
  random.seed,  
  alg_opt,  
  cdr3  
)
```

Arguments

`full_list` a list of clone lineages, represented as data.frames

csv	an indicator variable. TRUE if full_list argument is a list of csv files, FALSE otherwise
files	a list of data.frames. Each data.frame contains 2 columns, one that describes the sequences and the other which type of information (isotype or cluster) is considered in the analysis. All these cases are determined by the user.
distance_mat	a custom integer distance matrix, or NULL for using the default distance matrix (calculated based on the levenshtein distance, which counts the number of mutations between sequences).
clonal_frequency	a logical variable, TRUE if labeling of vertices is based on clonal frequency and FALSE otherwise.
scaleByClonalFreq	logical variable with TRUE if vertex size is scaled by the number of unique sequences per clone and FALSE otherwise.
weight	logical variable. When its value is FALSE, then the weights of outgoing edges from Germline node are set to 1. When its value is TRUE, the weights are set to the difference between the number of mutations among sequences in germline and connected nodes(value in the corresponding distance matrix) and the absolute value of the difference between the sequence lengths of germline and corresponding connected nodes. In both cases, weights of remaining edges are extracted from the distance matrix. Outgoing edges from Germline represent the number of mutations of sequences having as common ancestor the Germline.
tie_flag	a string, with options 'rand', 'full', 'close_to_germ', 'far_from_germ', 'close_path_to_germ', 'far_path_from_germ', 'most_expanded' and 'least_expanded' for removing edges when equal distance (tie) in distance matrix. 'rand' means random pruning in one of nodes, 'full' means keeping all nodes, 'close_to_germ' means pruning of node(s) farthest from germline (based on number of intermediate nodes), 'far_from_germ' means pruning of node(s) closest to germline (based on number of intermediate nodes), 'close_path_to_germ' means pruning of node(s) farthest from germline (based on edge path length), 'far_path_from_germ' means pruning of node(s) closest to germline (based on edge path length), 'most_expanded' means pruning of node(s) with the lowest B cell count(clonal frequency) and 'least_expanded', which means pruning of node(s) with the highest B cell count(clonal frequency). In cases of subsequent ties, a random node is selected.
scaleBybetweenness	logical variable with TRUE if vertex size is scaled by the vertex betweenness centrality.
scaleBycloseness_metr	logical variable with TRUE if vertex size is scaled by closeness centrality of vertices in graph.
opt	a string with options "isotype" and "cluster". The option "isotype" is utilized when the user desires to do an isotype analysis, while the selection of "cluster" denotes that an analysis based on transcriptome is requested.
random.seed	a random seed, specified by the user, when random sampling of sequences happens in each of the cases described in tie_flag argument.

alg_opt	a string denoting the version of the edge selection algorithm used in the construction of networks. Possible choices: "naive", "two-step".
cdr3	variable with values 0 if the user desires to select full length sequences (only when the input is a list of csv files), 1 for sequences in the CDR3 only (only when the input is a list of csv files) and NULL otherwise.

Value

graphs. A list of lists. E.g graphs[[1]][[1]] network: an igraph object, containing the first network in tree format. graphs[[1]][[2]] legend: contains the legend parameters of the first network. graphs[[1]][[3]] count.rand: contains the number of randomly considered nodes for the first network. graphs[[1]][[4]] adj.matrix: contains the adjacency matrix for the first network. graphs[[1]][[5]] distance.matrix: contains the distance matrix for the first network. graphs[[1]][[6]] cells.per.network: contains the number of cells for the first network. graphs[[1]][[7]] variants.per.network: contains the number of variants for the first network. graphs[[1]][[8]] variant.sequences: contains the sequences of the variants for the first network. graphs[[1]][[9]] cells.per.variant: contains the number of cells per variant (clonal frequency) for the first network. graphs[[1]][[10]] cell.indicies.per.variant: the indices of cells per variant for the first network. graphs[[1]][[11]] new.variant.names: contains the names of variants for the first network. graphs[[1]][[12]] germline.index: contains the index of germline sequence for the first network. graphs[[1]][[13]] isotype.per.variant: contains the isotypes corresponding to each variant for the first network. graphs[[1]][[14]] transcriptome.cluster.per.variant: contains the transcriptional clusters corresponding to each variant for the first network. graphs[[1]][[15]] isotype.per.cell: contains the isotype corresponding to each cell for the first network. graphs[[1]][[16]] transcriptome.cluster.per.cell: contains the transcriptional cluster corresponding to each cell for the first network.

See Also

ConvertStructure, CsvToDf, SubRepertoires, RemoveNets

Examples

```
## Not run:
AbForests_AntibodyForest(full_list = Platypus::new, csv=FALSE, files, clonal_frequency=TRUE,
scaleByClonalFreq=TRUE, weight=TRUE, tie_flag='close_to_germ',
scaleBybetweenness=FALSE, scaleBycloseness_metr=FALSE,
opt="cluster", alg_opt="0", cdr3=NULL)

## End(Not run)
```

AbForests_CompareForests

Comparison of distinct B cell repertoires

Description

CompareForests takes the output of AntibodyForest for 2 distinct repertoires and performs a comparison of these 2 repertoires.

Usage

```
AbForests_CompareForests(
  list1,
  list2,
  DAG,
  clonal_frequency,
  scaleByClonalFreq,
  weight,
  tie_flag,
  opt
)
```

Arguments

- | | |
|--------------------------------|--|
| <code>list1</code> | a list of lists. Each sublist contains an igraph object with the networks of the evolved B clonal lineages in tree format, their legend and the number of randomly considered nodes per network for Repertoire of 1 (Output of Antibody-Forest). E.g <code>list1[[1]][[1]]</code> is an igraph object, containing the first network of the evolved B clonal lineage in tree format. <code>list1[[1]][[2]]</code> contains the legend parameters of the first network of the evolved B clonal lineage. <code>list1[[1]][[3]]</code> is the number of randomly considered nodes for the first network of the evolved B clonal lineage. |
| <code>list2</code> | a list of lists. Each sublist contains an igraph object with the networks of the evolved B clonal lineages in tree format, their legend and the number of randomly considered nodes per network for Repertoire of 2 (Output of Antibody-Forest). E.g <code>list2[[1]][[1]]</code> is an igraph object, containing the first network of the evolved B clonal lineage in tree format. <code>list2[[1]][[2]]</code> contains the legend parameters of the first network of the evolved B clonal lineage. <code>list2[[1]][[3]]</code> is the number of randomly considered nodes for the first network of the evolved B clonal lineage. |
| <code>DAG</code> | a logical variable, when TRUE a directed acyclic graph is produced. |
| <code>clonal_frequency</code> | a logical variable, TRUE if labeling of vertices is based on clonal frequency and FALSE otherwise. |
| <code>scaleByClonalFreq</code> | logical variable with TRUE if vertex size is scaled by the number of unique sequences per clone and FALSE otherwise. |
| <code>weight</code> | logical variable. When its value is FALSE, then the weights of outgoing edges from Germline node are set to 1. When its value is TRUE, the weights are set to the difference between the number of mutations among sequences in germline and connected nodes(value in the corresponding distance matrix) and the absolute value of the difference between the sequence lengths of germline and corresponding connected nodes. In both cases, weights of remaining edges are extracted from the distance matrix. Outgoing edges from Germline represent the number of mutations of sequences having as common ancestor the Germline. |
| <code>tie_flag</code> | a string, with options 'rand', 'full', 'close_to_germ', 'far_from_germ', 'close_path_to_germ', 'far_path_from_germ', 'most_expanded' and 'least_expanded' for removing edges |

when equal distance (tie) in distance matrix. 'rand' means random pruning in one of nodes, 'full' means keeping all nodes, close_to_germ means pruning of node(s) farthest from germline (based on number of intermediate nodes), 'far_from_germ' means pruning of node(s) closest to germline (based on number of intermediate nodes), 'close_path_to_germ' means pruning of node(s) farthest from germline (based on edge path length), 'far_path_from_germ' means pruning of node(s) closest to germline (based on edge path length), 'most_expanded' means pruning of node(s) with the lowest B cell count (clonal frequency) and least_expanded, which means pruning of node(s) with the highest B cell count (clonal frequency). In cases of subsequent ties, a random node is selected.

opt a string with options "isotype" and "cluster". The option "isotype" is utilized when the user desires to do an isotype analysis, while the selection of "cluster" denotes that an analysis based on transcriptome is requested.

Value

combined_df. A data.frame that summarizes metrics for both repertoires. In particular, each row represents a single network and networks of both repertoires are combined row wise. Columns of combined_df are: Column1: Weighted.Longest.path.from.germline. Column2: Length.of.weighted.longest.shortest.path.from.germline. Column3: Unweighted.Longest.path.from.germline. Column4: Length.of.unweighted.longest.shortest.path.from.germline. Column5: Average.number.of.daughter.cells. Column6: Std.number.of.daughter.cells. Column7: Min.number.of.daughter.cells. Column8: Max.number.of.daughter.cells. Column9: Weighted.vertex.degree. Column10: Average.number.of.clusters/isotypes. Column11: Isotypes/Clusters.info. Column12: vertex.betweenness.centralitiy. Column13: edge.betweenness.centralitiy. Column14: closeness.centralitiy.of.vertices. Column15: global.clustering.coefficient. Column16: average.clustering.coefficient. Column17: Mean.clonal.expansion. If the labeling or scaling of nodes in graph is based on clonal frequency (arguments: clonal_frequency==TRUE or scaleByClonalFreq==TRUE), then combined_df contains also: Column18: Ratio.Number.of.edges.from.germline.to.each.node.with.clonal.frequency. Column19: Mean.Ratio.Number.of.edges.from.germline.to.each.node.with.clonal.frequency. Column20: Mean.number.of.edges.from.germline. Column21: Ratio.Total.path.length.from.germline.to.each.node.with.clonal.f. Column22: Mean.Ratio.Total.path.length.from.germline.to.each.node.with.clonal.frequency. Column23: Mean.Total.path.length.from.germline. Column24: Repertoire.id. Column25: Number.of.sequences.

isotype_info_rep1 A data.frame. It summarizes isotype/cluster info for repertoire 1.

isotype_info_rep2 A data.frame. It summarizes isotype/cluster info for repertoire 2.

See Also

AntibodyForest, ForestMetrics

Examples

```
## Not run:
CompareForests(list1,list2,DAG=TRUE,
clonal_frequency=TRUE,scaleByClonalFreq=TRUE,weight=TRUE,
tie_flag='close_to_germ',opt="cluster")

## End(Not run)
```

AbForests_ConvertStructure

Extract transcriptome/isotype information and B cell receptor sequences from single cell immune repertoire formatted as list of data.frames

Description

ConvertStructure alters a list of clone lineages, represented as data.frames and recovers the type of isotypes or transcriptional clusters and antibody sequences from these clone lineages. It can receive as input the original data or the output of SubRepertoiresByUniqueSeq or SubRepertoiresByCells. Then, the output list is used as input to the RemoveNets or AntibodyForest function.

Usage

```
AbForests_ConvertStructure(list, opt, cdr3)
```

Arguments

list	a list of data.frames. Each data.frame may contain information concerning full length heavy and light chain sequences, CDRH3 and CDRL3 sequences, the type of isotype and the transcriptional cluster that corresponds to each of these sequences.
opt	a string with options "isotype" and "cluster". The option "isotype" is utilized when the user desires to do an isotype analysis, while the selection of "cluster" denotes that an analysis based on transcriptome is requested.
cdr3	variable with values 0 if the user desires to select full length sequences (only when the input is a list of csv files), 1 for sequences in the CDR3 only (only when the input is a list of csv files) and NULL otherwise.

Value

list a list of data.frames. Each data.frame contains 2 columns, one that describes the sequences and the other which type of information (isotype or cluster) is considered in the analysis. All these cases are determined by the user.

See Also

RemoveNets, AntibodyForest

Examples

```
## Not run:  
ConvertStructure(list,opt="cluster",cdr3=NULL)  
ConvertStructure(list,opt="isotype",1)  
  
## End(Not run)
```

AbForests_CsvToDf *Convert list of csvs, to nested list of data.frames*

Description

CsvToDf converts a list of csv files, which are clone lineages to a list of data.frames.

Usage

```
AbForests_CsvToDf(files)
```

Arguments

files a list of csv files. Each csv file may contain information concerning full length heavy and light chain sequences, CDRH3 and CDRL3 sequences, the type of isotype and the transcriptional cluster that corresponds to each of these sequences.

Value

graphs a list of data.frames. Each data.frame contains 2 columns, one that describes the sequences and the other the type of information (isotype or cluster) is considered in the analysis. All these cases are determined by the user.

See Also

AntibodyForest

Examples

```
## Not run:  
CsvToDf(files)  
  
## End(Not run)
```

AbForests_ForestMetrics
Calculate metrics for networks

Description

ForestMetrics takes the output of AntibodyForest and calculates metrics for each of the networks.

Usage

```

AbForests_ForestMetrics(
  graphs,
  DAG,
  clonal_frequency,
  scaleByClonalFreq,
  weight,
  tie_flag,
  opt
)

```

Arguments

graphs	A list of lists. Each sublist contains an igraph object with the networks of the evolved B clonal lineages in tree format, their legend and the number of randomly considered nodes per network (Output of AntibodyForest function). E.g graphs[[1]][[1]] is an igraph object, containing the first network of the evolved B clonal lineage in tree format. graphs[[1]][[2]] contains the legend parameters of the first network of the evolved B clonal lineage. graphs[[1]][[3]] is the number of randomly considered nodes for the first network of the evolved B clonal lineage.
DAG	a logical variable, when TRUE a directed acyclic graph is produced.
clonal_frequency	a logical variable, TRUE if labeling of vertices is based on clonal frequency and FALSE otherwise.
scaleByClonalFreq	logical variable with TRUE if vertex size is scaled by the number of unique sequences per clone and FALSE otherwise.
weight	logical variable. When its value is FALSE, then the weights of outgoing edges from Germline node are set to 1. When its value is TRUE, the weights are set to the difference between the number of mutations among sequences in germline and connected nodes (value in the corresponding distance matrix) and the absolute value of the difference between the sequence lengths of germline and corresponding connected nodes. In both cases, weights of remaining edges are extracted from the distance matrix. Outgoing edges from Germline represent the number of mutations of sequences having as common ancestor the Germline.
tie_flag	a string, with options 'rand', 'full', 'close_to_germ', 'far_from_germ', 'close_path_to_germ', 'far_path_from_germ', 'most_expanded' and 'least_expanded' for removing edges when equal distance (tie) in distance matrix. 'rand' means random pruning in one of nodes, 'full' means keeping all nodes, close_to_germ means pruning of node(s) farthest from germline (based on number of intermediate nodes), 'far_from_germ' means pruning of node(s) closest to germline (based on number of intermediate nodes), 'close_path_to_germ' means pruning of node(s) farthest from germline (based on edge path length), 'far_path_from_germ' means pruning of node(s) closest to germline (based on edge path length), 'most_expanded' means pruning of node(s) with the lowest B cell count (clonal frequency) and least_expanded, which means pruning of node(s) with the highest B cell count (clonal frequency). In cases of subsequent ties, a random node is selected.

`opt` a string with options "isotype" and "cluster". The option "isotype" is utilized when the user desires to do an isotype analysis, while the selection of "cluster" denotes that an analysis based on transcriptome is requested.

Value

metrics. A list of lists. Each list contains various metrics for the quantification of networks. E.g `metrics[[1]][1]` is the weighted Longest path from germline for the first network. `metrics[[1]][2]` is the length of weighted longest shortest path from germline for the first network. `metrics[[1]][3]` is the unweighted Longest path from germline for the first network. `metrics[[1]][4]` is the length of unweighted longest shortest path from germline for the first network. `metrics[[1]][5]` is the weighted shortest path network for the first network. `metrics[[1]][6]` is the unweighted shortest path network for the first network. `metrics[[1]][7]` is the average number of daughter cells for the first network. `metrics[[1]][8]` is the std number of daughter cells for the first network. `metrics[[1]][9]` is the min number of daughter cells for the first network. `metrics[[1]][10]` is the max number of daughter cells for the first network. `metrics[[1]][11]` is a ggplot object that contains the plot of Degree Distribution of daughter cells for the first network. `metrics[[1]][12]` is the weighted vertex degree for the first network. `metrics[[1]][13]` is a ggplot object that contains the plot of unweighted Degree Distribution of daughter cells for the first network. `metrics[[1]][14]` is the average number of isotypes for the first network. `metrics[[1]][15]` is a ggplot object that contains the plot of Distribution of isotypes for the first network. `metrics[[1]][16]` is the Isotypes/Clusters info data.frame with columns Parent, Child and Parent-Child, which contains the type of isotypes/clusters for each pair of nodes (Parent-Child relationship in tree) found in the first network. `metrics[[1]][17]` is a ggplot object that contains the plot of Isotype/Cluster Directionality for the first network. In particular, the frequency of all types of isotypes/clusters for each pair of nodes in the tree is depicted. `metrics[[1]][18]` is the vertex betweenness centrality for the first network. It is defined by the number of geodesics (shortest paths) going through a vertex according to igraph documentation. `metrics[[1]][19]` is the edge betweenness centrality for the first network. It is defined by the number of geodesics (shortest paths) going through an edge according to igraph documentation. `metrics[[1]][20]` is the closeness centrality of vertices for the first network. The closeness centrality of a vertex is defined by the inverse of the average length of the shortest paths to/from all the other vertices in the graph according to igraph documentation. `metrics[[1]][21]` is a ggplot object that contains the plot of Path length from Germline vs Node Degree for the first network. `metrics[[1]][22]` is a ggplot object that contains the plot of Number of edges from Germline vs Node Degree for the first network. `metrics[[1]][23]` is an igraph object that contains the Isotype/Cluster transition network for the first network. `metrics[[1]][24]` is the global clustering coefficient for the first network. `metrics[[1]][25]` is the average clustering coefficient for the first network. `metrics[[1]][26]` is the mean clonal expansion for the first network, calculated as the mean of clonal frequencies of all vertices in the network. If the labeling or scaling of nodes in graph is based on clonal frequency (arguments: `clonal_frequency==TRUE` or `scaleBy-ClonalFreq==TRUE`), then `metrics[[1]][27]` is the ratio: Number of edges from germline to each node with clonal frequency for the first network. `metrics[[1]][28]` is the mean ratio: Number of edges from germline to each node with clonal frequency for the first network. `metrics[[1]][29]` is a ggplot object that contains the ratio of Number of edges from germline to each node with clonal frequency for the first network. `metrics[[1]][30]` is the mean number of edges from germline for the first network. `metrics[[1]][31]` is the ratio: Total path length from germline to each node with clonal frequency for the first network. `metrics[[1]][32]` is the mean ratio: Total path length from germline to each node with clonal frequency for the first network. `metrics[[1]][33]` is a ggplot object that contains the ratio of Total path length from germline to each node with clonal frequency

for the first network. `metrics[[1]][[34]]` is the mean Total path length from germline for the first network. `metrics[[2]][[1]]` is the weighted Longest path from germline for the second network.

See Also

AntibodyForest

Examples

```
## Not run:
ForestMetrics(graphs,DAG=TRUE,clonal_frequency=TRUE,scaleByClonalFreq=TRUE,
weight=TRUE,tie_flag='close_to_germ',opt="cluster")

## End(Not run)
```

AbForests_PlotGraphs *Plot igraph and ggplot objects*

Description

PlotGraphs takes as input the output of AntibodyForest or ForestMetrics functions and plots all corresponding networks in the single cell immune repertoire or the corresponding ggplot object, the user specifies, from all clone lineages. The function gives the option in the user to store each tree or ggplot object within the repertoire in pdf format.

Usage

```
AbForests_PlotGraphs(graphs, no_arg, topdf, filename)
```

Arguments

<code>graphs</code>	a list of networks (Output of AntibodyForest function) or metrics (Output of ForestMetrics function).
<code>no_arg</code>	element of list the user desires to plot : integer value,if the user desires to plot a metric and NULL, if the user desires to plot the networks.
<code>topdf</code>	logical value, TRUE if user wants to store plots in pdf format (the <code>no_arg</code> element of each list is saved in a separate page of the pdf).
<code>filename</code>	name of saved pdf file based on the user's preferences.

Value

Empty, output plots are written to file as specified by the user with the parameter filename

See Also

AntibodyForest, ForestMetrics

Examples

```
## Not run:  
PlotGraphs(graphs,no_arg=NULL,topdf=TRUE,filename)  
PlotGraphs(graphs,no_arg5,topdf=TRUE,filename)  
  
## End(Not run)
```

AbForests_PlyloToMatrix

Conversion of phylogenetic tree to distance matrix

Description

PlyloToMatrix converts a previously existing phylogenetic tree to a corresponding distance matrix using the cophenetic distance. Then, there is the option to utilize this custom distance matrix as an input distance matrix to AntibodyForest function. The user is responsible for specifying a correct and valid distance matrix. In particular, the size of distance matrix must match the number of sequences for each network in each repertoire.

Usage

```
AbForests_PlyloToMatrix(tree_name)
```

Arguments

tree_name a phylogenetic tree (phylo object).

Value

dist_mat The corresponding distance matrix uses the cophenetic distance between two observations that have been clustered. This distance is defined to be the intergroup dissimilarity at which the two observations are first combined into a single cluster.

See Also

AntibodyForest

Examples

```
## Not run:  
PlyloToMatrix(tree_name)  
  
## End(Not run)
```

AbForests_RemoveNets *Filter sub-repertoires with less than N unique sequences or with less than C unique cells*

Description

RemoveNets takes the output of SubRepertoires and performs the filtering of the 5 sub-repertoires. In particular, from these 5 sub-repertoires, networks with number of nodes or number of unique sequences below a specified threshold are eliminated.

Usage

```
AbForests_RemoveNets(
  list,
  opt,
  distance_mat,
  tie_flag,
  weight,
  N,
  C,
  random.seed,
  alg_opt
)
```

Arguments

list	a list of 5 sub-lists of data.frames. Each sub-list corresponds to the set of networks, in which a majority isotype is specified. list[[1]] or list\$list_IGHG contains the networks, in data.frame format, with more IGG isotypes, list[[2]] or list\$list_IGHA contains the networks, in data.frame format, with more IGA isotypes, list[[3]] or list\$list_IGHM contains the networks, in data.frame format, with more IGM isotypes, list[[4]] or list\$list_IGAG contains the networks, in data.frame format, with a tie in IGA and IGG isotypes and list[[5]] or list\$list_other contains the networks, in data.frame format, with other isotypes apart from the aforementioned combinations. In each sub-list, each data.frame represents a clone lineage and contains 2 columns, one that describes the antibody sequences and the other which type of information (isotype) is considered in the analysis. This list of data.frames has been generated by SubRepertoires function based on the initial data input and user's preferences.
opt	a string with options "isotype" and "cluster". The option "isotype" is utilized when the user desires to do an isotype analysis, while the selection of "cluster" denotes that an analysis based on transcriptome is requested.
distance_mat	a custom integer distance matrix, or NULL for using the default distance matrix (calculated based on the levenshtein distance, which counts the number of mutations between sequences). Given the phylogenetic tree, a custom-made distance matrix can be produced by PlyloToMatrix function.

<code>tie_flag</code>	a string, with options 'rand', 'full', 'close_to_germ', 'far_from_germ', 'close_path_to_germ', 'far_path_from_germ', 'most_expanded' and 'least_expanded' for removing edges when equal distance (tie) in distance matrix. 'rand' means random pruning in one of nodes, 'full' means keeping all nodes, 'close_to_germ' means pruning of node(s) farthest from germline (based on number of intermediate nodes), 'far_from_germ' means pruning of node(s) closest to germline (based on number of intermediate nodes), 'close_path_to_germ' means pruning of node(s) farthest from germline (based on edge path length), 'far_path_from_germ' means pruning of node(s) closest to germline (based on edge path length), 'most_expanded' means pruning of node(s) with the lowest B cell count (clonal frequency) and 'least_expanded', which means pruning of node(s) with the highest B cell count (clonal frequency). In cases of subsequent ties, a random node is selected.
<code>weight</code>	logical variable. When its value is FALSE, then the weights of outgoing edges from Germline node are set to 1. When its value is TRUE, the weights are set to the difference between the number of mutations among sequences in germline and connected nodes (value in the corresponding distance matrix) and the absolute value of the difference between the sequence lengths of germline and corresponding connected nodes. In both cases, weights of remaining edges are extracted from the distance matrix.
<code>N</code>	the threshold of unique sequences below which networks are removed.
<code>C</code>	the threshold of unique cells below which networks are removed.
<code>random.seed</code>	a random seed, specified by the user, when random sampling of sequences happens in each of the cases described in <code>tie_flag</code> argument.
<code>alg_opt</code>	a string denoting the version of the edge selection algorithm used in the construction of networks. "0" means the naive version and "1" the advanced one.

Value

list a nested list of 5 sub-lists of data.frames. Each sub-list corresponds to the reduced set of networks according to threshold N, in which a majority isotype is specified. `list[[1]]` or `list$list_IGHG` contains the networks, in data.frame format, with more IGG isotypes, `list[[2]]` or `list$list_IGHA` contains the networks, in data.frame format, with more IGA isotypes, `list[[3]]` or `list$list_IGHM` contains the networks, in data.frame format, with more IGM isotypes, `list[[4]]` or `list$list_IGAG` contains the networks, in data.frame format, with a tie in IGA and IGG isotypes and `list[[5]]` or `list$list_other` contains the networks, in data.frame format, with other isotypes apart from the aforementioned combinations.

See Also

SubRepertoires, SubRepertoiresByUniqueSeq, PlyloToMatrix, AntibodyForest

Examples

```
## Not run:
RemoveNets(list,opt="cluster",distance_mat=NULL,
tie_flag='close_to_germ',weight=TRUE,N=4,C=NULL,random.seed=165)

## End(Not run)
```

AbForests_SubRepertoiresByCells

Split single cell immune repertoire into sub-repertoires by isotype based on number of B cells

Description

SubRepertoiresByCells separates the single cell immune repertoire into 5 sub-repertoires taking into account the number of cells. The goal is to determine the majority isotype per each network in the immune repertoire. Therefore, each sub-repertoire is dominated by isotype IGG, IGA, IGM, other and if there is an equal number of IGA and IGG isotypes in a network, IGA-IGG category exists respectively. In particular, in case there exists a tie in the number of IGA and IGM, the network is considered to contain IGA as majority isotype, while the same number of IGG and IGM in the network categorize this network as containing IGG as majority isotype. The function receives the output of CsvToDf or original data and can be given as input to ConvertStructure function.

Usage

```
AbForests_SubRepertoiresByCells(list)
```

Arguments

`list` a list of data.frames. Each data.frame represents a clone lineage and separates initial input data into subsets of networks.

Value

list a nested list of 5 sub-lists of data.frames. Each sub-list corresponds to the set of networks, in which a majority isotype is specified. `list[[1]]` or `list$list_IGHG` contains the networks, in data.frame format, with more IGG isotypes, `list[[2]]` or `list$list_IGHA` contains the networks, in data.frame format, with more IGA isotypes, `list[[3]]` or `list$list_IGHM` contains the networks, in data.frame format, with more IGM isotypes, `list[[4]]` or `list$list_IGAG` contains the networks, in data.frame format, with a tie in IGA and IGG isotypes and `list[[5]]` or `list$list_other` contains the networks, in data.frame format, with other isotypes apart from the aforementioned combinations.

See Also

ConvertStructure, CsvToDf

Examples

```
## Not run:
SubRepertoiresByCells(list)

## End(Not run)
```

 AbForests_SubRepertoiresByUniqueSeq

Split single cell immune repertoire into sub-repertoires by isotype based on number of unique sequences

Description

SubRepertoiresByUniqueSeq separates the single cell immune repertoire into 5 sub-repertoires taking into account only unique sequences. The goal is to determine the majority isotype per each network in the immune repertoire. Therefore, each sub-repertoire is dominated by isotype IGG, IGA, IGM, other and if there is an equal number of IGA and IGG isotypes in a network, IGA-IGG category exists respectively. In particular, in case there exists a tie in the number of IGA and IGM, the network is considered to contain IGA as majority isotype, while the same number of IGG and IGM in the network categorize this network as containing IGG as majority isotype.

Usage

```
AbForests_SubRepertoiresByUniqueSeq(
  list,
  opt,
  distance_mat,
  tie_flag,
  weight,
  random.seed,
  alg_opt,
  cdr3
)
```

Arguments

<code>list</code>	a list of data.frames. Each data.frame represents a clone lineage and contains 2 columns, one that describes the antibody sequences and the other which type of information (isotype) is considered in the analysis. This list of data.frames has been generated by ConvertStructure function based on the initial data input or the output of CsvToDf and user's preferences.
<code>opt</code>	a string with options "isotype" and "cluster". The option "isotype" is utilized when the user desires to do an isotype analysis, while the selection of "cluster" denotes that an analysis based on transcriptome is requested.
<code>distance_mat</code>	a custom integer distance matrix, or NULL for using the default distance matrix (calculated based on the levenshtein distance, which counts the number of mutations between sequences). Given the phylogenetic tree, a custom-made distance matrix can be produced by PlyloToMatrix function.
<code>tie_flag</code>	a string, with options 'rand', 'full', 'close_to_germ', 'far_from_germ', 'close_path_to_germ', 'far_path_from_germ', 'most_expanded' and 'least_expanded' for removing edges when equal distance (tie) in distance matrix. 'rand' means random pruning

in one of nodes, 'full' means keeping all nodes, `close_to_germ` means pruning of node(s) farthest from germline (based on number of intermediate nodes), 'far_from_germ' means pruning of node(s) closest to germline (based on number of intermediate nodes), 'close_path_to_germ' means pruning of node(s) farthest from germline (based on edge path length), 'far_path_from_germ' means pruning of node(s) closest to germline (based on edge path length), 'most_expanded' means pruning of node(s) with the lowest B cell count (clonal frequency) and 'least_expanded', which means pruning of node(s) with the highest B cell count (clonal frequency). In cases of subsequent ties, a random node is selected.

<code>weight</code>	logical variable. When its value is FALSE, then the weights of outgoing edges from Germline node are set to 1. When its value is TRUE, the weights are set to the difference between the number of mutations among sequences in germline and connected nodes (value in the corresponding distance matrix) and the absolute value of the difference between the sequence lengths of germline and corresponding connected nodes. In both cases, weights of remaining edges are extracted from the distance matrix.
<code>random.seed</code>	a random seed, specified by the user, when random sampling of sequences happens in each of the cases described in <code>tie_flag</code> argument.
<code>alg_opt</code>	a string denoting the version of the edge selection algorithm used in the construction of networks. "0" means the naive version and "1" the advanced one.
<code>cdr3</code>	variable with values 0 if the user desires to select full length sequences (only when the input is a list of csv files), 1 for sequences in the CDR3 only (only when the input is a list of csv files) and NULL otherwise.

Value

list a nested list of 5 sub-lists of data.frames. Each sub-list corresponds to the set of networks, in which a majority isotype is specified. `list[[1]]` or `list$list_IGHG` contains the networks, in data.frame format, with more IGG isotypes, `list[[2]]` or `list$list_IGHA` contains the networks, in data.frame format, with more IGA isotypes, `list[[3]]` or `list$list_IGHM` contains the networks, in data.frame format, with more IGM isotypes, `list[[4]]` or `list$list_IGAG` contains the networks, in data.frame format, with a tie in IGA and IGG isotypes and `list[[5]]` or `list$list_other` contains the networks, in data.frame format, with other isotypes apart from the aforementioned combinations.

See Also

`AntibodyForest`, `ConvertStructure`, `CsvToDf`, `PlyloToMatrix`

Examples

```
## Not run:
SubRepertoiresByUniqueSeq(list,opt="isotype",distance_mat=NULL,
tie_flag='close_to_germ',weight=TRUE,random.seed=165,alg_opt="naive",cdr3=NULL)

## End(Not run)
```

 AbForests_UniqueAntibodyVariants

Count the number of unique antibody variants per clonal lineage

Description

UniqueAntibodyVariants calculates the number of unique antibody sequences, as dictated by the different grouping sequences strategy, for each network in the immune repertoire.

Usage

```
AbForests_UniqueAntibodyVariants(
  list,
  opt,
  distance_mat,
  tie_flag,
  weight,
  random.seed,
  alg_opt,
  cdr3
)
```

Arguments

<code>list</code>	a list of data.frames. Each data.frame represents a clone lineage and contains information on the antibody sequences and on the isotype/transcriptional cluster is considered in the analysis based the user's preferences.
<code>opt</code>	a string with options "isotype" and "cluster". The option "isotype" is utilized when the user desires to do an isotype analysis, while the selection of "cluster" denotes that an analysis based on transcriptome is requested.
<code>distance_mat</code>	a custom integer distance matrix, or NULL for using the default distance matrix (calculated based on the levenshtein distance, which counts the number of mutations between sequences). Given the phylogenetic tree, a custom-made distance matrix can be produced by PlyloToMatrix function.
<code>tie_flag</code>	a string, with options 'rand', 'full', 'close_to_germ', 'far_from_germ', 'close_path_to_germ', 'far_path_from_germ', 'most_expanded' and 'least_expanded' for removing edges when equal distance (tie) in distance matrix. 'rand' means random pruning in one of nodes, 'full' means keeping all nodes, 'close_to_germ' means pruning of node(s) farthest from germline (based on number of intermediate nodes), 'far_from_germ' means pruning of node(s) closest to germline (based on number of intermediate nodes), 'close_path_to_germ' means pruning of node(s) farthest from germline (based on edge path length), 'far_path_from_germ' means pruning of node(s) closest to germline (based on edge path length), 'most_expanded' means pruning of node(s) with the lowest B cell count(clonal frequency) and 'least_expanded', which means pruning of node(s) with the highest B cell count(clonal frequency). In cases of subsequent ties, a random node is selected.

weight	logical variable. When its value is FALSE, then the weights of outgoing edges from Germline node are set to 1. When its value is TRUE, the weights are set to the difference between the number of mutations among sequences in germline and connected nodes(value in the corresponding distance matrix) and the absolute value of the difference between the sequence lengths of germline and corresponding connected nodes. In both cases, weights of remaining edges are extracted from the distance matrix.
random.seed	a random seed, specified by the user, when random sampling of sequences happens in each of the cases described in tie_flag argument.
alg_opt	a string denoting the version of the edge selection algorithm used in the construction of networks. "0" means the naive version and "1" the advanced one.
cdr3	variable with values 0 if the user desires to select full length sequences (only when the input is a list of csv files), 1 for sequences in the CDR3 only (only when the input is a list of csv files) and NULL otherwise.

Value

uni_seq a vector, same size as list, which contains the number of unique antibody variants for each clonal lineage.

Examples

```
## Not run:
UniqueAntibodyVariants(list,opt="cluster",
distance_mat=NULL,tie_flag=close_to_germ,weight=TRUE,random.seed=165,alg_opt="naive",cdr3=NULL)

## End(Not run)
```

automate_GEX	<i>Automates the transcriptional analysis of the gene expression libraries from cellranger. This function will integrate multiple samples</i>
--------------	---

Description

Automates the transcriptional analysis of the gene expression libraries from cellranger. This function will integrate multiple samples

Usage

```
automate_GEX(
  GEX.outs.directory.list,
  GEX.list,
  integration.method,
  VDJ.gene.filter,
  mito.filter,
  norm.scale.factor,
  n.feature.rna,
```

```
n.count.rna.min,
n.count.rna.max,
n.variable.features,
cluster.resolution,
neighbor.dim,
mds.dim,
groups
)
```

Arguments

GEX.outs.directory.list

The path to the output of cellranger vdj runs. Multiple repertoires to be integrated together should be supplied as a character vector in the first element of a list. For example, if two separate VDJ repertoires should be integrated together (e.g. on the same tSNE plot), `GEX.outs.directory.list[[1]] <- c("my.VDJ1.path/outs/", "my.VDJ2.path/outs/...")` should be stored as input. If these repertoires should be analyzed separately, `>GEX.outs.directory.list[[1]] <- "my.VDJ1.path/outs/" >GEX.outs.directory.list[[2]] <- "my.VDJ2.path/outs/"` should be supplied.. This can be left blank if supplying the clonotypes and all_contig files directly as input. Multiple analyses can be stored

GEX.list

List containing the output from Seurat Read10x. This must be supplied if GEX.out.directory is not provided.

integration.method

String specifying which data normalization and integration pipeline should be used. Default is "scale.data", which correspondings to the ScaleData function internal to harmony package. 'sct' specifies SCTransform from the Seurat package. "harmony" should be specified to perform harmony integration. This method requires the harmony package from bioconductor.

VDJ.gene.filter

Logical indicating if variable genes from the b cell receptor and t cell receptor should be removed from the analysis. True is highly recommended to avoid clonal families clustering together.

mito.filter

Numeric specifying which percent of genes are allowed to be composed of mitochondrial genes. This value may require visual inspection and can be specific to each sequencing experiment. Users can visualize the percentage of genes corresponding to mitochondrial genes using the function "investigate_mitochondial_genes".

norm.scale.factor

Scaling factor for the standard Seurat pipeline. Default is set to 10000 as reported in Seurat documentation.

n.feature.rna

Numeric that specifies which cells should be filtered out due to low number of detected genes. Default is set to 0. Seurat standard pipeline uses 2000.

n.count.rna.min

Numeric that specifies which cells should be filtered out due to low RNA count.Default is set to 0. Seurat standard pipeline without VDJ information uses 200.

n.count.rna.max

Numeric that specifies which cells should be filtered out due to high RNA count.Default is set to infinity. Seurat standard pipeline without VDJ information uses 2500.

n.variable.features	Numeric specifying the number of variable features. Default set to 2000 as specified in Seurat standard pipeline.
cluster.resolution	Numeric specifying the resolution that will be supplied to Seurat's FindClusters function. Default is set to 0.5. Increasing this number will increase the number of distinct Seurat clusters. Suggested to examine multiple parameters to ensure gene signatures differentiating clusters remains constant.
neighbor.dim	Numeric vector specifying which dimensions should be supplied in the Find-Neighbors function from Seurat. Default input is '1:10'.
mds.dim	Numeric vector specifying which dimensions should be supplied into dimensional reduction techniques in Seurat and Harmony. Default input is '1:10'.
groups	Integer specifying the groups of the different samples. This is needed if there are multiple biological replicates for a given condition sequenced and aligned through cellranger separately.

Value

Returns a processed Seurat object containing transcriptional information from all samples which can be supplied to the VDJ_GEX_integrate function.

Examples

```
## Not run:
automate_GEX(out_directory=fullRepertoire.output,
  rep.size=3*length(unlist(fullRepertoire.output[[1]])),
  distribution="identical",
  with.germline="FALSE")

## End(Not run)
```

Bcell_sequences_example_tree
Example csv file 1

Description

Example csv file 1

Usage

```
Bcell_sequences_example_tree
```

Format

An object of class data.frame with 170 rows and 1 columns.

References

R package Platypus : <https://doi.org/10.1093/nargab/lqab023>

Bcell_tree_2	<i>Example csv file 2</i>
--------------	---------------------------

Description

Example csv file 2

Usage

```
Bcell_tree_2
```

Format

An object of class data.frame with 85 rows and 1 columns.

References

R package Platypus:<https://doi.org/10.1093/nargab/lqab023>

call_MIXCR	<i>Extracts information on the VDJRegion level using MiXCR. This function assumes the user can run an executable instance of MiXCR and is eligible to use MiXCR as determined by license agreements. The VDJRegion corresponds to the recombined heavy and light chain loci starting from framework region 1 (FR1) and extending to frame work region 4 (FR4). This can be useful for extracting full-length sequences ready to clone and further calculating somatic hypermutation occurrences.</i>
------------	--

Description

Extracts information on the VDJRegion level using MiXCR. This function assumes the user can run an executable instance of MiXCR and is eligible to use MiXCR as determined by license agreements. The VDJRegion corresponds to the recombined heavy and light chain loci starting from framework region 1 (FR1) and extending to frame work region 4 (FR4). This can be useful for extracting full-length sequences ready to clone and further calculating somatic hypermutation occurrences.

Usage

```
call_MIXCR(VDJ.per.clone, mixcr.directory, species)
```

Arguments

VDJ.per.clone	The output from the VDJ_per_clone function. This object should have information regarding the contigs and clonotype_ids for each cell.
mixcr.directory	The directory containing an executable version of MiXCR. This must be downloaded separately and is under a separate license.
species	Either "mmu" for mouse or "hsa" for human. These use the default germline genes for both species contained in MIXCR.

Value

Returns a nested list containing VDJRegion information as determined by MIXCR. The outer list corresponds to the individual repertoires in the same structure as the input VDJ.per.clone. The inner list corresponds to each clonal family, as determined by either the VDJ_clonotype function or the default nucleotide clonotyping produced by cellranger. Each element in the inner list corresponds to a dataframe containing repertoire information such as isotype, CDR sequences, mean number of UMIs. This output can be supplied to further package functions such as VDJ_extract_sequences and VDJ_GEX_integrate.

See Also

VDJ_extract_sequences

Examples

```
## Not run:
call_MIXCR(VDJ.per.clone = VDJ.per.clone.output
,mixcr.directory = "~/Downloads/mixcr-3.0.12/mixcr",species = "mmu")

## End(Not run)
```

class_switch_prob_hum *class_switch_prob_hum* The probability matrix of class switching for human b cells. The row names of the matrix are the isotypes the cell is switching from, the column names are the isotypes the cell is switching to. All B cells start from IGHM, and switch to one of the other isotypes or remain the same.

Description

class_switch_prob_hum The probability matrix of class switching for human b cells. The row names of the matrix are the isotypes the cell is switching from, the column names are the isotypes the cell is switching to. All B cells start from IGHM, and switch to one of the other isotypes or remain the same.

Usage

```
data("class_switch_prob_hum")
```

Format

A 8*8 matrix. The row and column names are "IGHM", "IGHD", "IGHG1", "IGHG2", "IGHG3", "IGHG4", "IGHE", "IGHA". The probability for a cell to switch from "IGHM" to "IGHD" is the value at class_switch_prob_hum[1,2].

class_switch_prob_mus *class_switch_prob_mus* The probability matrix of class switching for mouse b cells. The row names of the matrix are the isotypes the cell is switching from, the column names are the isotypes the cell is switching to. All B cells start from IGHM, and switch to one of the other isotypes or remain the same.

Description

class_switch_prob_mus The probability matrix of class switching for mouse b cells. The row names of the matrix are the isotypes the cell is switching from, the column names are the isotypes the cell is switching to. All B cells start from IGHM, and switch to one of the other isotypes or remain the same.

Usage

```
data("class_switch_prob_mus")
```

Format

A 9*9 matrix. The row and column names are "IGHM", "IGHD", "IGHG1", "IGHG2A", "IGHG2B", "IGHG2C", "IGHG3", "IGHG4". The probability for a cell to switch from "IGHM" to "IGHD" is the value at class_switch_prob_mus[1,2].

clonofreq *Plot clonal frequency barplot of the output simulated data*

Description

Plot the top abundant clonal frequencies in a barplot with ggplot2

Usage

```
clonofreq(clonotypes, top.n, y.limit)
```

Arguments

clonotypes	The clonotypes dataframe, which is the second element in the simulation output list.
top.n	The top n abundant clones to be shown in the plot. If missing, all clones will be shown.
y.limit	The upper limit for y axis in the plot.

Value

top abundant clonal frequencies in a barplot with ggplot2

clonofreq.isotype.data

Get information about the clonotype counts grouped by isotype.

Description

Return

Usage

```
clonofreq.isotype.data(all.contig.annotations, top.n)
```

Arguments

all.contig.annotations

The output dataframe all_contig_annotation from function simulate.repertoire.

top.n

The top n abundant clones to be shown in the plot. If missing, all clones will be shown.

Value

dataframes containing the top n abundant clonotypes and their frequency and isotype information for further processing.

clonofreq.isotype.plot

Get information about the clonotype counts grouped by isotype.

Description

Plot a stacked barplot for clonotype counts grouped by isotype.

Usage

```
clonofreq.isotype.plot(all.contig.annotations, top.n, y.limit, colors)
```

Arguments

all.contig.annotations	The output dataframe all_contig_annotation from function simulate.repertoire.
top.n	The top n abundant clones to be shown in the plot. If missing, all clones will be shown.
y.limit	The upper limit for y axis in the plot.
colors	A named character vector of colors, the names are the isotypes. If missing, the default has 11 colors corresponding to the default isotype names.

Value

a stacked barplot for clonotype counts grouped by isotype

clonofreq.trans.data *Get information about the clonotype counts grouped by transcriptome state(cell type).*

Description

Dataframe with clonotype counts grouped by transcriptome state(cell type).

Usage

```
clonofreq.trans.data(all.contig.annotations, history, trans.names, top.n)
```

Arguments

all.contig.annotations	The output dataframe all_contig_annotation from function simulate.repertoire.
history	The dataframe history from simulate output.
trans.names	The names of cell types which are used in transcriptome.switch.prob argument in the simulation.
top.n	The top n abundant clones to be shown in the plot. If missing, all clones will be shown.

Value

a dataframe with clonotype counts grouped by transcriptome state(cell type).

`clonofreq.trans.plot` *Get information about the clonotype counts grouped by transcriptome state(cell type).*

Description

Plot a stacked barplot for clonotype counts grouped by transcriptome state(cell type).

Usage

```
clonofreq.trans.plot(
  all.contig.annotations,
  history,
  trans.names,
  top.n,
  y.limit,
  colors
)
```

Arguments

<code>all.contig.annotations</code>	The output dataframe <code>all_contig_annotation</code> from function <code>simulate.repertoire</code> .
<code>history</code>	The dataframe <code>history</code> from <code>simulate</code> output.
<code>trans.names</code>	The names of cell types which are used in <code>transcriptome.switch.prob</code> argument in the simulation.
<code>top.n</code>	The top n abundant clones to be shown in the plot. If missing, all clones will be shown.
<code>y.limit</code>	The upper limit for y axis in the plot.
<code>colors</code>	A named character vector of colors, the names are the isotypes. If missing, the default has 11 colors corresponding to the default isotype names.

Value

a stacked barplot for clonotype counts grouped by transcriptome state(cell type).

`cluster.id.igraph` *Get clone network igrhps colored by seurat cluster id.*

Description

Get clone network igrhps colored by seurat cluster id.

Usage

```
cluster.id.igraph(meta.data, history, igraph.index, empty.node)
```

Arguments

meta.data	the meta.data dataframe from the Seurat object of the simulation. The object should be pre-processed and has cluster ids in the meta.data.
history	The dataframe 'history' from the simulation output.
igraph.index	The list 'igraph.index' from the simulation output.
empty.node	If TRUE, there will be empty node in igraph. if FALSE, the empty node will be deleted.

Value

a list of clone network igraphs colored by seurat cluster id

colors	<i>colors</i> A vector of characters specifying colors used in igraph phylogenetic tree. Default colors: "#66C2A5", "#FC8D62", "#8DA0CB", "#E78AC3", "#A6D854"
--------	--

Description

colors A vector of characters specifying colors used in igraph phylogenetic tree. Default colors: "#66C2A5", "#FC8D62", "#8DA0CB", "#E78AC3", "#A6D854"

Usage

```
data("colors")
```

Format

a character vector

Echidna_vae_generate *Simulate B or T cell receptor sequences by variational autoencodes(VAEs) trained with experimental data.*

Description

Simulate B or T cell receptor sequences by variational autoencodes(VAEs) trained with experimental data.

Usage

```
Echidna_vae_generate(
    sequence,
    n.train,
    n.sample,
    batch.size,
    latent.dim,
    intermediate.dim,
    epochs,
    epsilon.std,
    null.threshold
)
```

Arguments

sequence	a vector of seuqnece the model to be trained on
n.train	number of sequence to be used in training set, the rest will be in testing set
n.sample	number of new sequence to generate from VAE model
batch.size	set to larger to save time, set to smaller to same computing power
latent.dim	parameter used in VAE model
intermediate.dim	parameter used in VAE model
epochs	parameter used in VAE model
epsilon.std	parameter used in VAE model
null.threshold	threshold of predicted value to be considered as an existing base, default is 0.05. When generated sequence is too short, lower this threshold.

Value

A simulated VDJ repertoire on the basis of the input experimental repertoire

get.avr.mut.data	<i>Get information about somatic hypermutation in the simulation. This function return a barplot showing the average mutation.</i>
------------------	--

Description

Get information about somatic hypermutation in the simulation. This function return a barplot showing the average mutation.

Usage

```
get.avr.mut.data(igraph.index.attr, history, clonotype.select, level)
```

Arguments

igraph.index.attr	A list "igraph.index.attr" from the simulation output.
history	A dataframe "history" from the simulation output.
clonotype.select	The selected clonotype index, can be the output of the function "select.top".
level	Can be "clone" or "cell". If "clone", the function will return average mutation on unique variant level. Otherwise it will return on cell level.

Value

a bar plot showing the average mutation on clone or cell level.

get.avr.mut.plot	<i>Get information about somatic hypermutation in the simulation. This function return a barplot showing the average mutation.</i>
------------------	--

Description

Get information about somatic hypermutation in the simulation. This function return a barplot showing the average mutation.

Usage

```
get.avr.mut.plot(igraph.index.attr, history, clonotype.select, level, y.limit)
```

Arguments

<code>igraph.index.attr</code>	A list "igraph.index.attr" from the simulation output.
<code>history</code>	A dataframe "history" from the simulation output.
<code>clonotype.select</code>	The selected clonotype index, can be the output of the function "select.top".
<code>level</code>	Can be "node" or "cell". If "node", the function will return average mutation on unique variant level. Otherwise it will return on cell level.
<code>y.limit</code>	The upper limit for y axis in the plot.

Value

a barplot showing the average mutation per node (same heavy and light chain set) or per cell.

`get.barplot.errorbar` *Return a barplot of mean and standard error bar of certain value of each clone.*

Description

Return a barplot of mean and standard error bar of certain value of each clone.

Usage

```
get.barplot.errorbar(data, y.lab, y.limit)
```

Arguments

<code>data</code>	A dataframe. Columns are different simulations, rows are the top clones. The first row is the top abundant clone.
<code>y.lab</code>	A string specifies the y lable name of the barplot.
<code>y.limit</code>	The upper limit for y axis in the plot.

Value

a barplot of mean and standard error bar of certain value of each clone.

get.elbow	<i>Get the seurat object from simulated transcriptome output.</i>
-----------	---

Description

Get the seurat object from simulated transcriptome output.

Usage

```
get.elbow(data)
```

Arguments

data The output "transcriptome" dataframe from simulation output.

Value

the seurat object from simulated transcriptome output.

get.n.node.data	<i>Get the number of unique variants in each clone in a vector. The output is the vector representing the numbers of unique variants.</i>
-----------------	---

Description

Get the number of unique variants in each clone in a vector. The output is the vector representing the numbers of unique variants.

Usage

```
get.n.node.data(data, clonotype.select)
```

Arguments

data The output "igraph.index.attr" list from simulation output.
clonotype.select The index of the clones to be shown. If missing, all clones will be included.

Value

the number of unique variants in each clone in a vector. The output is the vector representing the numbers of unique variants.

<code>get.n.node.plot</code>	<i>Get the number of unique variants in each clone in a vector and the barplot. The first item in the output is the vector representing the numbers of unique variants, the second item is the barplot.</i>
------------------------------	---

Description

Get the number of unique variants in each clone in a vector and the barplot. The first item in the output is the vector representing the numbers of unique variants, the second item is the barplot.

Usage

```
get.n.node.plot(igraph.index.attr, clonotype.select, y.limit)
```

Arguments

<code>igraph.index.attr</code>	The output "igraph.index.attr" list from simulation output.
<code>clonotype.select</code>	The index of the clones to be shown. If missing, all clones will be included.
<code>y.limit</code>	The upper limit for y axis in the plot.

Value

the number of unique variants in each clone in a vector and the barplot. The first item in the output is the vector representing the numbers of unique variants, the second item is the barplot.

<code>get.seq.distance</code>	<i>Computing sequence distance according to the number of unmatched bases.</i>
-------------------------------	--

Description

Computing sequence distance according to the number of unmatched bases.

Usage

```
get.seq.distance(germline, sequence)
```

Arguments

<code>germline</code>	A string representing the germline sequence.
<code>sequence</code>	A string of the sequence to be compared, which has the same length as germline.

Value

the number of unmatched bases in 2 sequences.

get.umap	<i>Further process the seurat object from simulated transcriptome output and make UMAP ready for plotting.</i>
----------	--

Description

Further process the seurat object from simulated transcriptome output and make UMAP ready for plotting.

Usage

```
get.umap(gex, d, reso)
```

Arguments

gex	output from get.elbow function.
d	dims argument of in Seurat::FindNeighbors() and Seurat::RunUMAP
reso	resolution argument in Seurat::FindClusters()

Value

Further processed seurat object from simulated transcriptome output with UMAP ready for plotting.

get.vgu.matrix	<i>Get paired v gene heavy chain and light chain matrix on clonotype level. A v gene usage pheatmap can be obtain by p<-pheatmap::pheatmap(vgu_matrix,show_colnames= T, main = "V Gene Usage"), where the vgu_matrix is the output of this function.</i>
----------------	---

Description

Get paired v gene heavy chain and light chain matrix on clonotype level. A v gene usage pheatmap can be obtain by p<-pheatmap::pheatmap(vgu_matrix,show_colnames= T, main = "V Gene Usage"), where the vgu_matrix is the output of this function.

Usage

```
get.vgu.matrix(all.contig.annotations, level)
```

Arguments

all.contig.annotations	The dataframe "all_contig_annotation" from simulation output.
level	Can be "clone" or "cell". If "clone", the function will return paired v gene usage matrix on clonotype level. Otherwise it will return on cell level.

Value

a paired v gene heavy chain and light chain matrix on clonotype level.

GEX_automate	<i>Automates the transcriptional analysis of the gene expression libraries from cellranger. This function will integrate multiple samples</i>
--------------	---

Description

Automates the transcriptional analysis of the gene expression libraries from cellranger. This function will integrate multiple samples

Usage

```
GEX_automate(
  GEX.outs.directory.list,
  GEX.list,
  integration.method,
  VDJ.gene.filter,
  mito.filter,
  norm.scale.factor,
  n.feature.rna,
  n.count.rna.min,
  n.count.rna.max,
  n.variable.features,
  cluster.resolution,
  neighbor.dim,
  mds.dim,
  groups
)
```

Arguments

GEX.outs.directory.list

The path to the output of cellranger vdj runs. Multiple repertoires to be integrated together should be supplied as a character vector in the first element of a list. For example, if two separate VDJ repertoires should be integrated together (e.g. on the same tSNE plot), `GEX.outs.directory.list[[1]] <- c("my.VDJ1.path/outs/", "my.VDJ2.path/outs/...")` should be stored as input. If these repertoires should be analyzed separately, `>GEX.outs.directory.list[[1]] <- "my.VDJ1.path/outs/" >GEX.outs.directory.list[[2]] <- "my.VDJ2.path/outs/"` should be supplied. This can be left blank if supplying the clonotypes and all_contig files directly as input. Multiple analyses can be stored

GEX.list

List containing the output from Seurat Read10x. This must be supplied if GEX.out.directory is not provided.

<code>integration.method</code>	String specifying which data normalization and integration pipeline should be used. Default is "scale.data", which corresponds to the ScaleData function internal to harmony package. 'sct' specifies SCTransform from the Seurat package. "harmony" should be specified to perform harmony integration. This method requires the harmony package from bioconductor.
<code>VDJ.gene.filter</code>	Logical indicating if variable genes from the b cell receptor and t cell receptor should be removed from the analysis. True is highly recommended to avoid clonal families clustering together.
<code>mito.filter</code>	Numeric specifying which percent of genes are allowed to be composed of mitochondrial genes. This value may require visual inspection and can be specific to each sequencing experiment. Users can visualize the percentage of genes corresponding to mitochondrial genes using the function "investigate_mitochondrial_genes".
<code>norm.scale.factor</code>	Scaling factor for the standard Seurat pipeline. Default is set to 10000 as reported in Seurat documentation.
<code>n.feature.rna</code>	Numeric that specifies which cells should be filtered out due to low number of detected genes. Default is set to 0. Seurat standard pipeline uses 2000.
<code>n.count.rna.min</code>	Numeric that specifies which cells should be filtered out due to low RNA count. Default is set to 0. Seurat standard pipeline without VDJ information uses 200.
<code>n.count.rna.max</code>	Numeric that specifies which cells should be filtered out due to high RNA count. Default is set to infinity. Seurat standard pipeline without VDJ information uses 2500.
<code>n.variable.features</code>	Numeric specifying the number of variable features. Default set to 2000 as specified in Seurat standard pipeline.
<code>cluster.resolution</code>	Numeric specifying the resolution that will be supplied to Seurat's FindClusters function. Default is set to 0.5. Increasing this number will increase the number of distinct Seurat clusters. Suggested to examine multiple parameters to ensure gene signatures differentiating clusters remains constant.
<code>neighbor.dim</code>	Numeric vector specifying which dimensions should be supplied in the FindNeighbors function from Seurat. Default input is '1:10'.
<code>mds.dim</code>	Numeric vector specifying which dimensions should be supplied into dimensional reduction techniques in Seurat and Harmony. Default input is '1:10'.
<code>groups</code>	Integer specifying the groups of the different samples. This is needed if there are multiple biological replicates for a given condition sequenced and aligned through cellranger separately.

Value

Returns a processed Seurat object containing transcriptional information from all samples which can be supplied to the VDJ_GEX_integrate function.

Examples

```
## Not run:
automate_GEX(out_directory=fullRepertoire.output,
  rep.size=3*length(unlist(fullRepertoire.output[[1]])),
  distribution="identical",
  with.germline="FALSE")

## End(Not run)
```

GEX_clonotype	<i>Platypus V2: Integrates VDJ and gene expression libraries by providing cluster membership seq_per_vdj object and the index of the cell in the Seurat RNA-seq object.</i>
---------------	---

Description

Platypus V2: Integrates VDJ and gene expression libraries by providing cluster membership seq_per_vdj object and the index of the cell in the Seurat RNA-seq object.

Usage

```
GEX_clonotype(GEX.object, VDJ.per.clone)
```

Arguments

GEX.object	A single seurat object from automate_GEX function. This will likely be supplied as automate_GEX.output[[1]].
VDJ.per.clone	Output from the VDJ_per_clone function. Each element in the list should be found in the output from the automate_GEX function.

Value

Returns a dataframe containing repertoire information, such as isotype, CDR sequences, mean number of UMIs. This output can be supplied to further packages VDJ_extract_sequences and VDJ_GEX_integrate

Examples

```
## Not run:
GEX_clonotype(GEX.object=automate.GEX.output[[1]], VDJ.per.clone=vdj.per.clone.output)

## End(Not run)
```

GEX_cluster_genes	<i>Extracts the differentially expressed genes between two samples. This function uses the FindMarkers function from the Seurat package. Further parameter control can be accomplished by calling the function directly on the output of automate_GEX or VDJ_GEX_matrix</i>
-------------------	---

Description

Extracts the differentially expressed genes between two samples. This function uses the FindMarkers function from the Seurat package. Further parameter control can be accomplished by calling the function directly on the output of automate_GEX or VDJ_GEX_matrix

Usage

```
GEX_cluster_genes(GEX, min.pct, filter, base, platypus.version)
```

Arguments

GEX	Output Seurat object of either automate_GEX for platypus.version v2 or VDJ_GEX_matrix for platypus.version v3 (usually VDJ_GEX_matrix.output[[2]])
min.pct	The minimum percentage of cells expressing a gene in either of the two groups to be compared. Default is 0.25
filter	Character vector of initials of the genes to be filtered. Default is c("MT-", "RPL", "RPS"), which filters mitochondrial and ribosomal genes.
base	The base with respect to which logarithms are computed. Default: 2
platypus.version	is set automatically

Value

Returns a dataframe containing the output from the FindMarkers function, which contains information regarding the genes that are differentially regulated, statistics (p value and log fold change), and the percent of cells expressing the particular gene. Each element in the list corresponds to the clusters in numerical order. For example, the first element in the list output[[1]] corresponds to the genes differentially expressed in cluster 0 in GEX

Examples

```
#Platypus version v2
#GEX_cluster_genes(GEX = automate_GEX_output[[i]], min.pct = .25
#, filter = c("MT-", "RPL", "RPS"))

#Platypus version v3
GEX_cluster_genes(GEX = subset(Platypus::small_vgm[[2]], seurat_clusters %in% c(0,1)), min.pct = .25
, filter = c("MT-", "RPL", "RPS"))
```

GEX_cluster_genes_heatmap

Produces a heatmap displaying the expression of the top genes that define each cluster in the Seurat object. The output heatmap is derived from DoHeatmap from Seurat and thereby can be edited using typical ggplot interactions. The number of genes per cluster and the number of cells to display can be specified by the user. Either the log fold change or the p value can be used to select the top n genes.

Description

Produces a heatmap displaying the expression of the top genes that define each cluster in the Seurat object. The output heatmap is derived from DoHeatmap from Seurat and thereby can be edited using typical ggplot interactions. The number of genes per cluster and the number of cells to display can be specified by the user. Either the log fold change or the p value can be used to select the top n genes.

Usage

```
GEX_cluster_genes_heatmap(  
  GEX,  
  GEX_cluster_genes.output,  
  n.genes.per.cluster,  
  metric,  
  max.cell,  
  group.colors,  
  slot,  
  platypus.version  
)
```

Arguments

- | | |
|--------------------------|--|
| GEX | Output Seurat object of either automate_GEX for platypus.version v2 or of VDJ_GEX_matrix for platypus.version v3 (usually VDJ_GEX_matrix.output[[2]]) |
| GEX_cluster_genes.output | The output from the GEX_cluster_genes function - this should be a list with each list element corresponding to the genes, p values, logFC, pct expression for the genes differentially regulated for each cluster. |
| n.genes.per.cluster | An integer value determining how many genes per cluster to display in the output heatmap. This number should be adjusted based on the number of clusters. Too many genes per cluster and clusters may cause a problem with the heatmap function in Seurat. |
| metric | The metric that dictates which are the top n genes returned. Possible options are "p.value" (default), "avg_logFC", "top_logFC", "bottom_logFC". "top_logFC" returns the top expressed genes for each cluster, whereas "bottom_logFC" returns the least expressed genes per cluster-both by log fold change. |

max.cell	The max number of cells to display in the heatmap for each cluster, which corresponds to the number of columns. Default is set to 100 cells per cluster.
group.colors	Optional character vector. Array of colors with the same length as GEX_cluster_genes.output to color bars above the heatmap. Defaults to rainbow palette
slot	Seurat object slot from which to plot gene expression data.
platypus.version	is set automatically

Value

Returns a heatmap from the function DoHeatmap from the package Seurat, which is a ggplot object that can be modified or plotted. The number of genes is determined by the n.genes parameter and the number of cells per cluster is determined by the max.cell argument. This function gives a visual description of the top genes differentially expressed in each cluster.

Examples

```
## Not run:
#For Platypus version 2
cluster_defining_gene_heatmap <- GEX_cluster_genes_heatmap(GEX = automate_GEX_output[[i]]
,GEX_cluster_genes.output=GEX_cluster_genes_output
,n.genes.per.cluster=5,metric="p.value",max.cell=5)

#For Platypus version 3

cluster_defining_gene_heatmap <- GEX_cluster_genes_heatmap(GEX = VDJ_GEX_matrix.output[[2]]
,GEX_cluster_genes.output=GEX_cluster_genes_output
,n.genes.per.cluster=5,metric="p.value",max.cell=5)

## End(Not run)
```

GEX_cluster_membership

Plots the cluster membership for each of the distinct samples in the Seurat object from the automate_GEX function. The distinct samples are determined by "sample_id" field in the Seurat object.

Description

Plots the cluster membership for each of the distinct samples in the Seurat object from the automate_GEX function. The distinct samples are determined by "sample_id" field in the Seurat object.

Usage

```
GEX_cluster_membership(GEX, by.group, platypus.version)
```

Arguments

GEX	Output Seurat object containing gene expression data from <code>automate_GEX</code> (<code>platypus.version = "v2"</code>) or <code>VDJ_GEX_matrix</code> (<code>platypus.version = "v3"</code> , usually <code>VDJ_GEX_matrix.output[[2]]</code>) that contained at least two distinct biological samples. The different biological samples correspond to integer values (<code>v2</code>) or factor values (<code>v3</code>) in the order of the working directories initially supplied to the <code>automate_GEX</code> function.
<code>by.group</code>	Logical indicating whether to look at the cluster distribution per group (using the <code>group_id</code> column). Default is set to <code>FALSE</code> .
<code>platypus.version</code>	Version of <code>platypus</code> to use. Defaults to <code>"v2"</code> . If an output of the <code>GEX_automate</code> function is supplied, set to <code>"v2"</code> . If an output of the <code>VDJ_GEX_matrix</code> function is supplied set to <code>"v3"</code>

Value

Returns a `ggplot` object in which the values on the x axis correspond to each cluster found in the Seurat object. The y axis corresponds to the percentage of cells found in each cluster. The bar and color corresponds to the distinct `sample_id`.

Examples

```
#Platypus v2
#GEX_cluster_membership(GEX=automate_GEX_out[[2]], platypus.version = "v2")
#Platypus v3
GEX_cluster_membership(GEX= Platypus::small_vgm[[2]], platypus.version = "v3")
```

GEX_coexpression_coefficient

Returns either a plot or numeric data of coexpression levels of selected genes. Coexpression % is calculated as the quotient of double positive cells (counts > 0) and the sum of total cells positive for either genes.

Description

Returns either a plot or numeric data of coexpression levels of selected genes. Coexpression % is calculated as the quotient of double positive cells (counts > 0) and the sum of total cells positive for either genes.

Usage

```
GEX_coexpression_coefficient(GEX, genes, subsample.n, plot.dotmap)
```

Arguments

GEX	GEX seurat object generated with VDJ_GEX_matrix (VDJ_GEX_matrix.output[[2]])
genes	Character vector. At least 2 genes present in rownames(GEX). Use "all" to include all genes. The number of comparisons to make is the length(genes)! (factorial). More than 100 genes are not recommended.
subsample.n	Integer. Number of cells to subsample. If set to 100, 100 cells will be randomly sampled for the calculation
plot.dotmap	Boolean. Whether to return a plot

Value

Returns a dataframe if plot.dotmap == F or a ggplot if plot.dotmap == T detailing the coexpression levels of selected genes within the given cell population

Examples

```
#To return a dataframe with coefficients
#GEX_coexpression_coefficient(GEX = VDJ_GEX_matrix.output[[2]]
#, genes = c("CD19", "EBF1", "SDC1"), subsample.n = "none", plot.dotmap = FALSE)

#To return a dotplot detailing coexpression and overall expression
GEX_coexpression_coefficient(GEX = Platypus::small_vgm[[2]]
, genes = c("CD19", "CD83"), subsample.n = "none", plot.dotmap = FALSE)
```

GEX_DEgenes	<i>Extracts the differentially expressed genes between two groups of cells. These groups are defined as cells having either of two entries (group1, group2) in the grouping.column of the input Seurat object metadata This function uses the FindMarkers function from the Seurat package.</i>
-------------	---

Description

Extracts the differentially expressed genes between two groups of cells. These groups are defined as cells having either of two entries (group1, group2) in the grouping.column of the input Seurat object metadata This function uses the FindMarkers function from the Seurat package.

Usage

```
GEX_DEgenes(
  GEX,
  FindMarkers.out,
  grouping.column,
  group1,
  group2,
  min.pct,
  filter,
```

```

return.plot,
logFC,
color.p.threshold,
color.log.threshold,
color.by.threshold,
up.genes,
down.genes,
base,
label.n.top.genes,
genes.to.label
)

```

Arguments

GEX	Output Seurat object from <code>automate_GEX</code> or <code>VDJ_GEX_matrix_function</code> (<code>VDJ_GEX_matrix.output[[2]]</code>) function that contained at least two distinct biological groups.
FindMarkers.out	OPTIONAL: the output of the <code>FindMarkers</code> function. This skips the DEG calculation step and outputs desired plots. All plotting parameters function as normal. Grouping parameters and <code>min.pct</code> are ignored.
grouping.column	Character. A column name of <code>GEX@meta.data</code> . In this column, <code>group1</code> and <code>group2</code> should be found. Defaults to <code>"sample_id"</code> . Could also be set to <code>"seurat_clusters"</code> to generate DEGs between cells of 2 chosen clusters.
group1	either character or integer specifying the first group of cells that should be compared. (e.g. <code>"s1"</code> if <code>sample_id</code> is used as <code>grouping.column</code>)
group2	either character or integer specifying the first group of cells that should be compared. (e.g. <code>"s2"</code> if <code>sample_id</code> is used as <code>grouping.column</code>)
min.pct	The minimum percentage of cells expressing a gene in either of the two groups to be compared.
filter	Character vector of initials of the genes to be filtered. Default is <code>c("MT-", "RPL", "RPS")</code> , which filters mitochondrial and ribosomal genes.
return.plot	Character specifying if a <code>"heatmap"</code> , or a <code>"volcano"</code> or <code>"none"</code> is to be returned. If not <code>"none"</code> then <code>@return</code> is a list where the first element is a dataframe and the second a plot (see <code>@return</code>). Defaults to <code>none</code>
logFC	Logical specifying whether the genes will be displayed based on <code>logFC</code> (TRUE) or <code>pvalue</code> (FALSE).
color.p.threshold	numeric specifying the adjusted p-value threshold for <code>geom_points</code> to be colored. Default is set to 0.01.
color.log.threshold	numeric specifying the absolute <code>logFC</code> threshold for <code>geom_points</code> to be colored. Default is set to 0.25.
color.by.threshold	Boolean. Set to TRUE to color by <code>color.p.threshold</code> and <code>color.log.threshold</code> . Set to FALSE for a continuous color scale by fold change.

up.genes	FOR HEATMAP Integer specifying the number of upregulated genes to be shown.
down.genes	FOR HEATMAP Integer specifying the number of downregulated genes to be shown.
base	The base with respect to which logarithms are computed. Default: 2
label.n.top.genes	FOR VOLCANO Integer. How many top genes to label either by Fold change (if logFC == TRUE) or by p.value (if logFC == FALSE). More than 50 are not recommended. Also works in conjunction with genes.to.label
genes.to.label	FOR VOLCANO Character vector of genes to label irregardless of their p value.

Value

Returns a list with out[[1]] = a dataframe containing the output from the FindMarkers function, which contains information regarding the genes that are differentially regulated, statistics (p value and log fold change), and the percent of cells expressing the particular gene for both groups. out[[2]] = either a "heatmap" (set return.plot accordingly), or "volcano" plot

Examples

```
#Basic run between two samples
GEX_DEgenes(GEX = Platypus::small_vgm[[2]],min.pct = .25,
group1 = "s1",group2 = "s2")

#Getting DEGs between two seurat clusters
#GEX_DEgenes(GEX = Platypus::small_vgm[[2]],min.pct = .25,
#grouping.column = "seurat_clusters",group1 = "0",group2 = "1")

#Plotting a heatmap by foldchange of sample markers
#GEX_DEgenes(GEX = VDJ_GEX_matrix.output[[2]]
#,min.pct = .25,group1 = "s1",group2 = "s2", return.plot = "heatmap"
#, up.genes = 10, down.genes = 10, logFC = TRUE)

#Plotting volcano by p value of sample markers. Label additional genes of interest
#GEX_DEgenes(GEX = VDJ_GEX_matrix.output[[2]],min.pct = .25
#,group1 = "s1",group2 = "s2", return.plot = "volcano", logFC = FALSE
#, label.n.top.genes = 40, genes.to.label = c("CD28", "ICOS"))

#Generate a heatmap from an already existing FindMarkers output
#GEX_DEgenes(GEX = VDJ_GEX_matrix.output[[2]]
#, FindMarkers.out = FindMarkers.output.dataframe, return.plot = "heatmap"
#, up.genes = 10, down.genes = 10, logFC = TRUE, platypus.version = "v3")
```

GEX_DEgenes_persample *!Only for Platypus version v2. For more flexibility and platypus v3 please refer to GEX_Degenes. Extracts the differentially expressed genes between two samples. This function uses the FindMarkers function from the Seurat package. Further parameter control can be accomplished by calling the function directly on the output of automate_GEX and further extracting sample information from the "sample_id" component of the Seurat object.*

Description

!Only for Platypus version v2. For more flexibility and platypus v3 please refer to GEX_Degenes. Extracts the differentially expressed genes between two samples. This function uses the FindMarkers function from the Seurat package. Further parameter control can be accomplished by calling the function directly on the output of automate_GEX and further extracting sample information from the "sample_id" component of the Seurat object.

Usage

```
GEX_DEgenes_persample(  
  automate.GEX,  
  min.pct,  
  sample1,  
  sample2,  
  by.group,  
  filter,  
  return.plot,  
  logFC,  
  up.genes,  
  down.genes,  
  base  
)
```

Arguments

automate.GEX	Output Seurat object from automate_GEX function that contained at least two distinct biological samples. The differential biological samples correspond to integer values in the order of the working directories initially supplied to the automate_GEX function.
min.pct	The minimum percentage of cells expressing a gene in either of the two groups to be compared.
sample1	either character or integer specifying the first sample that should be compared.
sample2	either character or integer specifying the first sample that should be compared.
by.group	Logical specifying if groups should be used instead of samples. If TRUE, then the argument in sample1 and sample2 will correspond to cells found in the groups from sample1 or sample2.

filter	Character vector of initials of the genes to be filtered. Default is c("MT-", "RPL", "RPS"), which filters mitochondrial and ribosomal genes.
return.plot	Logical specifying if a heatmap of the DEX genes is to be returned. If TRUE then @return is a list where the first element is a dataframe and the second a heatmap (see @return)
logFC	Logical specifying whether the genes will be displayed based on logFC (TRUE) or pvalue (FALSE).
up.genes	Integer specifying the number of upregulated genes to be shown.
down.genes	Integer specifying the number of downregulated genes to be shown.
base	The base with respect to which logarithms are computed. Default: 2

Value

Returns a dataframe containing the output from the FindMarkers function, which contains information regarding the genes that are differentially regulated, statistics (p value and log fold change), and the percent of cells expressing the particular gene for both groups.

Examples

```
## Not run:
GEX_DEgenes_persample(automate.GEX=automate.GEX.output[[i]]
,min.pct = .25,sample1 = "1",sample2 = "2")

## End(Not run)
```

GEX_dottile_plot	<i>Outputs a dotplot for gene expression, where the color of each dot is scaled by the gene expression level and the size is scaled by the % of cells positive for the gene</i>
------------------	---

Description

Outputs a dotplot for gene expression, where the color of each dot is scaled by the gene expression level and the size is scaled by the % of cells positive for the gene

Usage

```
GEX_dottile_plot(GEX, genes, group.by, threshold.to.plot, platypus.version)
```

Arguments

GEX	GEX seurat object generated with VDJ_GEX_matrix
genes	Character vector. Genes of those in rownames(GEX) to plot. Can be any number, but more than 30 is discouraged because of cluttering
group.by	Character. Name of a column in GEX@meta.data to split the plot by. If set to "none", a plot with a single column will be produced.

threshold.to.plot

Integer 1-100. % of cells which must be expressing the feature to plot a point.
If below, the field will be left empty

platypus.version

This is coded for \"v3\" only, but in practice any Seurat Object can be fed in

Value

Returns a ggplot object where the dot size indicates the percentage of expressing cells and the dot color indicates the expression level.

Examples

```
#To return a plot detailing the expression of common genes by seurat cluster
GEX_dottile_plot(GEX = Platypus::small_vgm[[2]], genes = c("CD19", "CD83"),
group.by = "seurat_clusters", threshold.to.plot = 5)
```

GEX_GOterm

Runs a GO term analysis on a submitted list of genes. Works with the output of GEX_topN_DE_genes_per_cluster or a custom list of genes to obtain GOterms.

Description

Runs a GO term analysis on a submitted list of genes. Works with the output of GEX_topN_DE_genes_per_cluster or a custom list of genes to obtain GOterms.

Usage

```
GEX_GOterm(
  GEX.cluster.genes.output,
  topNgenes,
  ontology,
  species,
  up.or.down,
  MT.Rb.filter,
  kegg,
  go.plots,
  top.N.go.terms.plots,
  kegg.plots,
  top.N.kegg.terms.plots
)
```

Arguments

<code>GEX.cluster.genes.output</code>	Either output of <code>Platypus::GEX_cluster_genes</code> or custom character vector containing gene symbols. A custom gene list will not be further filtered or ordered.
<code>topNgenes</code>	How many of the most significant up or down regulated genes should be considered for GO term analysis. All genes will be used if left empty.
<code>ontology</code>	Ontology used for the GO terms. "MF", "BP" or "CC" possible. Default: "BP"
<code>species</code>	The species the genes belong to. Default: "Mm".
<code>up.or.down</code>	Whether up or downregulated genes should be used for GO term analysis if <code>GEX_cluster_genes</code> output is used. Default: "up"
<code>MT.Rb.filter</code>	logical, if mitochondrial and ribosomal genes should be filtered out.
<code>kegg</code>	logical, if KEGG pathway analysis should be conducted. Requires internet connection. Default: False.
<code>go.plots</code>	logical, if top GO-terms should be visualized. Default: False. If True, for each cluster the top N (<code>top.N.GO.terms.plots</code>) Go-terms for each cluster will be plotted to the working directory and saved as a list element. Plots are made both based on <code>padj</code> and <code>ratio</code> .
<code>top.N.go.terms.plots</code>	The number of most significant GO-terms to be included in the <code>go.plots</code> . Default: 10.
<code>kegg.plots</code>	logical, if top KEGG-terms should be visualized. Default: False. If True, for each cluster the top N (<code>top.N.kegg.terms.plots</code>) KEGG-terms for each cluster will be plotted to the working directory and saved as a list element. Plots are made both based on <code>padj</code> and <code>ratio</code> .
<code>top.N.kegg.terms.plots</code>	The number of most significant KEGG-terms to be included in the <code>kegg.plots</code> . Default: 10.

Value

Returns a list of data frames and plots containing the TopGO and the TopKEGG output containing the significant GO/KEGG terms and their visualizations.

Examples

```
## Not run:

GEX_GOterm(DE_genes_cluster,MT.Rb.filter = TRUE, species= "Mm", ontology = "MF")
GEX_GOterm(rownames(DE_genes_cluster[[1]]),MT.Rb.filter = TRUE, species= "Mm",
ontology = "BP", go.plots = TRUE)

#Install the needed database with
#if (!requireNamespace("BiocManager", quietly = TRUE))
#install.packages("BiocManager")
#BiocManager::install("org.Mm.eg.db")

## End(Not run)
```

GEX_GSEA	<i>Conducts a Gene Set Enrichment Analysis (GSEA) on a set of genes submitted in a data frame with a metric each. Works with the output of GEX_genes_cluster or a custom data frame containing the gene symbols either in a column "symbols" or as rownames and a metric for each gene. The name of the column containing the metric has to be declared via the input metric.colname.</i>
----------	---

Description

Conducts a Gene Set Enrichment Analysis (GSEA) on a set of genes submitted in a data frame with a metric each. Works with the output of GEX_genes_cluster or a custom data frame containing the gene symbols either in a column "symbols" or as rownames and a metric for each gene. The name of the column containing the metric has to be declared via the input metric.colname.

Usage

```
GEX_GSEA(
  GEX.cluster.genes.output,
  MT.Rb.filter,
  filter,
  path.to.pathways,
  metric.colname,
  pval.adj.cutoff,
  Enrichment.Plots,
  my.own.geneset,
  eps,
  platypus.version,
  verbose
)
```

Arguments

GEX.cluster.genes.output	Data frame containing the list of gene symbols and a metric. Function works directly with GEX_cluster_genes output.
MT.Rb.filter	Logical, should Mitotic and Ribosomal genes be filtered out of the geneset. True by default.
filter	Character vector containing the identifying symbol sequence for the genes which should be filtered out, if MT.Rb.filter == T. By default set to c("MT-", "RPL", "RPS").
path.to.pathways	Either a path to gmt file containing the gene sets (can be downloaded from MSigDB) or vector where first element specifies species and second element specifies the MSigDB collection abbreviation. E.g.: c("Homo sapiens", "H"). Mouse C7 (immunologic signature) gene set will be used by default.

<code>metric.colname</code>	Name of column which contains the metric used for the ranking of the submitted genelist. "avg_logFC" is used by default.
<code>pval.adj.cutoff</code>	Only genes with a more significant adjusted pvalue are considered. Default: 0.001
<code>Enrichment.Plots</code>	List of Gene-set names which should be plotted as Enrichment plots in addition to the top 10 Up and Downregulated Genesets.
<code>my.own.geneset</code>	A list, where each element contains a gene list and is named with the corresponding pathway name. Default is set to FALSE, so that gene sets from MSigDB are used. Should not contain ".gmt" in name.
<code>eps</code>	Numeric, specifying boundary for calculating the p value in the GSEA.
<code>platypus.version</code>	Function works with V2 and V3, no need to set this parameter.
<code>verbose</code>	Print run parameters and status to console

Value

Returns a list containing a tibble with the gene sets and their enrichment scores and Enrichment plots. List element `[[1]]`: Dataframe with Genesets and statistics. `[[2]]`: Enrichment plots of top10 Up regulated genesets. `[[3]]`: Enrichment plots of top10 Down regulated genesets. `[[4]]`: Enrichment plots of submitted gene-sets in parameter `Enrichment.Plot`.

Examples

```
## Not run:
df <- GEX_cluster_genes(gex_combined[[1]])

#Using gmt file to perform gsea
output <- GEX_GSEA(GEX.cluster.genes.output = df[[1]], MT.Rb.filter = TRUE
, path.to.pathways = "./c5.go.bp.v7.2.symbols.gmt")
cowplot::plot_grid(plotlist=output[[2]], ncol=2)
View(gex_gsea[[1]])

#Directly downloading gene set collection from MSigDB to perform gsea
output <- GEX_GSEA(GEX.cluster.genes.output = df[[1]], MT.Rb.filter = TRUE
, path.to.pathways = c("Mus musculus", "C7"))

#Using your own gene list to perform gsea
output <- GEX_GSEA(GEX.cluster.genes.output = df[[1]], MT.Rb.filter = TRUE
, my.own.geneset = my_geneset)

## End(Not run)
```

GEX_heatmap	<i>Produces a heatmap containing gene expression information at the clonotype level. The rows correspond to different genes that can either be determined by pre-made sets of B or T cell markers, or can be customized by the user. The columns correspond to individual cells and the colors correspond to the different clonotype families.</i>
-------------	--

Description

Produces a heatmap containing gene expression information at the clonotype level. The rows correspond to different genes that can either be determined by pre-made sets of B or T cell markers, or can be customized by the user. The columns correspond to individual cells and the colors correspond to the different clonotype families.

Usage

```
GEX_heatmap(
  GEX,
  b.or.t,
  sample.index,
  clone.rank.threshold,
  custom.array,
  slot
)
```

Arguments

GEX	A single <code>seurat</code> object from <code>clonotype_GEX</code> function corresponding to all of the samples in a single <code>VDJ_analyze</code> object. This will likely be supplied as <code>clonotype_GEX.output[[i]]</code> if there were multiple, distinct transcriptomes.
<code>b.or.t</code>	Logical indicating if B or T cell gene panel should be used.
<code>sample.index</code>	Corresponds to which repertoire should be used in the case that the length of <code>clonotype.list</code> has a length greater than 1. The transcriptional profiles from only one repertoire can be plotted at a time.
<code>clone.rank.threshold</code>	A numeric that specifies the threshold clonal rank that specifies which clonotypes to extract transcriptome information from. For example, if 10 is supplied then the gene expression for the top ten clones included on the heatmap, separated by clonotype.
<code>custom.array</code>	Corresponds to which repertoire should be used in the case that the length of <code>clonotype.list</code> has a length greater than 1. The transcriptional profiles from only one repertoire can be plotted at a time.
<code>slot</code>	<code>Seurat</code> data slot from which to plot values. Can be "raw.data", "data" or "scale.data"

Value

Returns a heatmap via `Seurat::DoHeatmap` of gene expression per clonotype

See Also

VDJ_extract_sequences

Examples

```
#prep the small_vgm sample dataset
small_vgm <- Platypus::small_vgm
small_vgm[[2]]$clone_rank <- c(1:nrow(small_vgm[[2]]@meta.data))
GEX_heatmap(GEX = small_vgm[[2]],b.or.t = "custom"
,clone.rank.threshold = 1,sample.index = "s1"
,custom.array = c("CD24A","CD83"), slot = "data")
```

GEX_pairwise_DEGs	<i>Produces and saves a list of volcano plots with each showing differentially expressed genes between pairs groups. If e.g. <code>seurat_clusters</code> used as <code>group.by</code>, a plot will be generated for every pairwise comparison of clusters. For large numbers of this may take longer to run. Only available for platypus v3</i>
-------------------	---

Description

Produces and saves a list of volcano plots with each showing differentially expressed genes between pairs groups. If e.g. `seurat_clusters` used as `group.by`, a plot will be generated for every pairwise comparison of clusters. For large numbers of this may take longer to run. Only available for platypus v3

Usage

```
GEX_pairwise_DEGs(
  GEX,
  group.by,
  min.pct,
  RP.MT.filter,
  label.n.top.genes,
  genes.to.label,
  save.plot
)
```

Arguments

GEX	Output Seurat object of the <code>VDJ_GEX_matrix</code> function (<code>VDJ_GEX_matrix.output[[2]]</code>)
group.by	Character. Defaults to "seurat_clusters" Column name of <code>GEX@meta.data</code> to use for pairwise comparisons. More than 20 groups are discouraged.
min.pct	Numeric. Defaults to 0.25 passed to <code>Seurat::FindMarkers</code>
RP.MT.filter	Boolean. Defaults to True. If True, mitochondrial and ribosomal genes are filtered out from the output of <code>Seurat::FindMarkers</code>

<code>label.n.top.genes</code>	Integer. Defaults to 50. Defines how many genes are labelled via <code>geom_text_repel</code> . Genes are ordered by adjusted p value and the first <code>label.n.genes</code> are labelled
<code>genes.to.label</code>	Character vector. Defaults to "none". Vector of gene names to plot independently of their p value. Can be used in combination with <code>label.n.genes</code> .
<code>save.plot</code>	Boolean. Defaults to True. Whether to save plots as appropriately named .png files

Value

A nested list with `out[[i]][[1]]` being ggplot volcano plots and `out[[i]][[2]]` being source DEG dataframes.

Examples

```
GEX_pairwise_DEGs(GEX = Platypus::small_vgm[[2]],group.by = "sample_id"
,min.pct = 0.25,RP.MT.filter = TRUE,label.n.top.genes = 2,genes.to.label = c("CD24A")
,save.plot = FALSE)
```

<code>GEX_phenotype</code>	<i>Integrates VDJ and gene expression libraries by providing cluster membership <code>seq_per_vdj</code> object and the index of the cell in the Seurat RNA-seq object.</i>
----------------------------	---

Description

Integrates VDJ and gene expression libraries by providing cluster membership `seq_per_vdj` object and the index of the cell in the Seurat RNA-seq object.

Usage

```
GEX_phenotype(seurat.object, cell.state.names, cell.state.markers, default)
```

Arguments

<code>seurat.object</code>	A single seurat object from <code>automate_GEX</code> function
<code>cell.state.names</code>	Character vector containing the gene names for each state. ; is used to use multiple markers within a single gene state. Different vector elements correspond to different states. Order must match <code>cell.state.names</code> containing the <code>c("CD4+;CD44-","CD4+;IL7R+;CD44+")</code> .
<code>cell.state.markers</code>	Character vector containing the cell state labels defined by the markers in <code>cell.state.markers</code> parameter. Example is <code>c("NaiveCd4","MemoryCd4")</code> .
<code>default</code>	Default is TRUE - will use premade gene sets and cell states.

Value

Returns a stacked barplot that visualizes the seurat cluster membership for different cell phenotypes.

Examples

```
GEX_phenotype.test <- GEX_phenotype(seurat.object = Platypus::small_vgm[[2]])
```

GEX_phenotype_per_clone

Integrates VDJ and gene expression libraries by providing cluster membership seq_per_vdj object and the index of the cell in the Seurat RNA-seq object. ! For platypus.version == "v3" and VDJ_GEX_matrix output the function will iterate over entries in the sample_id column of the GEX by default.

Description

Integrates VDJ and gene expression libraries by providing cluster membership seq_per_vdj object and the index of the cell in the Seurat RNA-seq object. ! For platypus.version == "v3" and VDJ_GEX_matrix output the function will iterate over entries in the sample_id column of the GEX by default.

Usage

```
GEX_phenotype_per_clone(
  GEX,
  clonotype.ids,
  global.clonotypes,
  GEX.group.by,
  GEX.clonotypes,
  platypus.version
)
```

Arguments

GEX	For platypus.version == "v3" the GEX object from the output of the VDJ_GEX_matrix function (VDJ_GEX_matrix.output $\setminus[2\setminus]$). For platypus.version == "v2" a single seurat object from automate_GEX function after labeling cell phenotypes using the GEX_phenotype function.
clonotype.ids	For platypus.version == "v2" Output from either VDJ_analyze or VDJ_clonotype functions. This list should correspond to a single GEX.list object, in which each list element in clonotype.list is found in the GEX.object. Furthermore, these repertoires should be found in the automate_GEX library.
global.clonotypes	Boolean. Defaults to FALSE. Set to True if clonotyping has been done across samples

GEX.group.by For platypus.version == "v3". Character. Column name of the GEX@meta.data to group barplot by. Defaults to `seurat_clusters`

GEX.clonotypes For platypus.version == "v3". Numeric vector with ids of clonotypes to plot e.g. `c(1,2,3,4)`. Can also be set to "topclones"

platypus.version Set to either "v2" or "v3" depending on whether supplying GEX_automate or VDJ_GEX_matrix\[\[2\]\] objects. Defaults to "v3"

Value

Returns a stacked barplot that visualizes the seurat cluster membership for different cell phenotypes.

Examples

```
#For testing: only a single clonotype in two samples
small_vgm_cl <- Platypus::small_vgm
small_vgm_cl[[2]]$clonotype_id_10x <- "clonotype1"
GEX_phenotype_per_clone(GEX = small_vgm_cl[[2]]
, GEX.clonotypes = c(1), GEX.group.by = "seurat_clusters", platypus.version = "v3")
```

GEX_proportions_barplot

Plots proportions of a group of cells within a secondary group of cells. E.g. The proportions of samples in seurat clusters, or the proportions of samples in defined cell subtypes

Description

Plots proportions of a group of cells within a secondary group of cells. E.g. The proportions of samples in seurat clusters, or the proportions of samples in defined cell subtypes

Usage

```
GEX_proportions_barplot(GEX, source.group, target.group, stacked.plot, verbose)
```

Arguments

GEX GEX Seurat object generated with VDJ_GEX_matrix (VDJ_GEX_matrix.output[[2]])

source.group Character. A column name of the GEX@meta.data with the group of which proportions should be plotted

target.group Character. A column name of the GEX@meta.data with the group to calculate proportions within. If unsure, see examples for clarification

stacked.plot Boolean. Defaults to FALSE. Whether to return a stacked barplot, with the y axis representing the % of cells of the target group. If set to FALSE a normal barplot (position = "dodge") will be returned with the y axis representing the % of cells of the source group

verbose Print information about factor levels and ordering to console

Value

Returns a ggplot barplot showing cell proportions by source and target group.

Examples

```
#To return a normal barplot which shows the % of cells of
#each sample contained in each cluster
GEX_proportions_barplot(GEX = Platypus::small_vgm[[2]], source.group = "sample_id"
, target.group = "seurat_clusters",stacked.plot = FALSE)
```

```
#To return a stacked barplot which shows the % of cells of each
#cluster attributed to each sample
GEX_proportions_barplot(GEX = Platypus::small_vgm[[2]],
source.group = "sample_id", target.group = "seurat_clusters"
,stacked.plot = TRUE)
```

GEX_scatter_coexpression

Clonal frequency plot displaying clonal expansion for either T and B cells with Platypus v3 input.

Description

Clonal frequency plot displaying clonal expansion for either T and B cells with Platypus v3 input.

Usage

```
GEX_scatter_coexpression(GEX, gene.1, gene.2, color.theme)
```

Arguments

GEX	GEX seurat object generated with VDJ_GEX_matrix
gene.1	Character. Name of a gene in rownames(VDJ.matrix)
gene.2	Character. Name of a gene in rownames(VDJ.matrix)
color.theme	Character. A color to use for the composite plot

Value

Returns a gridplot showing coexpression scatterplot as well as histograms of gene.1 and gene.2

Examples

```
gene1 <- "CD24A"
gene2 <- "CD83"
GEX_scatter_coexpression(GEX = Platypus::small_vgm[[2]], gene1,gene2)
```

GEX_topN_DE_genes_per_cluster

Organizes the top N genes that define each Seurat cluster and converts them into a single dataframe. This can be useful for obtaining insight into cluster-specific phenotypes.

Description

Organizes the top N genes that define each Seurat cluster and converts them into a single dataframe. This can be useful for obtaining insight into cluster-specific phenotypes.

Usage

```
GEX_topN_DE_genes_per_cluster(GEX_cluster_genes.output, n.genes, by_FC, filter)
```

Arguments

GEX_cluster_genes.output	The output from the GEX_cluster_genes function - this should be a list with each list element corresponding to the genes, p values, logFC, pct expression for the genes differentially regulated for each cluster.
n.genes	The number of genes to be selected from each cluster. If n.genes is higher than the number of cells in a cluster then it is silently adjusted to be
by_FC	Logical indicating if the top n genes are selected based on the logFC value instead of p value. Default is FALSE.
filter	Character vector of initials of the genes to be filtered. Default is c("MT-", "RPL", "RPS"), which filters mitochondrial and ribosomal genes.

Value

Returns a dataframe in which the top N genes defining each cluster based on differential expression are selected.

Examples

```
## Not run:
GEX_topDE_genes_per_cluster(GEX_cluster_genes.output=list_of_genes_per_cluster
, n.genes=20, by_FC=FALSE, filter=c("MT-", "RPS", "RPL"))

## End(Not run)
```

GEX_visualize_clones *!Only for platypus version v2. For platypus v3 refer to: VDJ_GEX_overlay_clones() Visualize selected clonotypes on the tSNE or UMAP projection.*

Description

!Only for platypus version v2. For platypus v3 refer to: VDJ_GEX_overlay_clones() Visualize selected clonotypes on the tSNE or UMAP projection.

Usage

```
GEX_visualize_clones(  
  GEX.list,  
  VDJ.GEX.integrate.list,  
  highlight.type,  
  highlight.number,  
  reduction  
)
```

Arguments

`GEX.list` list of Seurat objects, output of the `automate_GEX` function.

`VDJ.GEX.integrate.list` Output of the `VDJ_GEX_integrate` function.

`highlight.type` (Optional) either "None" if representation highlighted by cluster, "clonotype" if want to highlight most expanded clonotypes, or "sample" if several samples are within the same Seurat object. Default is None.

`highlight.number` (Optional) an integer or list of integers representing the number of most expanded clonotypes or samples one wants to select eg 4 to highlight the 4th most expanded clonotype or 2:5 to highlight the top 2 to top 5 most expanded clonotype. Only compatible with `highlight.type` "clonotype" or "sample", will be ignored if type is "None". Default is 1.

`reduction` (Optional) Reduction to plot, either "tsne", "umap", or "harmony". Default is "tsne".

Value

concatenated ggplot2 plot with selected clonotypes highlighted (if None, the coloring is according to the clustering).

Examples

```
## Not run:
GEX_visualize_clones(GEX.list=automate_GEX.output,
  VDJ.per.clone=VDJ_per_clone.output,
  highlight.type="clonotype",
  highlight.number=1:4,
  reduction="umap")

## End(Not run)
```

GEX_volcano	<i>Plots a volcano plot from the output of the FindMarkers function from the Seurat package or the GEX_cluster_genes function alternatively.</i>
-------------	--

Description

Plots a volcano plot from the output of the FindMarkers function from the Seurat package or the GEX_cluster_genes function alternatively.

Usage

```
GEX_volcano(
  DEGs.input,
  input.type,
  condition.1,
  condition.2,
  explicit.title,
  RP.MT.filter,
  color.p.threshold,
  color.log.threshold,
  label.p.threshold,
  label.logfc.threshold,
  n.label.up,
  n.label.down,
  by.logFC,
  maximum.overlaps,
  plot.adj.pvalue
)
```

Arguments

DEGs.input	Either output data frame from the FindMarkers function from the Seurat package or GEX_cluster_genes list output.
input.type	Character specifying the input type as either "findmarkers" or "cluster.genes". Defaults to "cluster.genes"

<code>condition.1</code>	either character or integer specifying <code>ident.1</code> that was used in the <code>FindMarkers</code> function from the <code>Seurat</code> package. Should be left empty when using the <code>GEX_cluster_genes</code> output.
<code>condition.2</code>	either character or integer specifying <code>ident.2</code> that was used in the <code>FindMarkers</code> function from the <code>Seurat</code> package. Should be left empty when using the <code>GEX_cluster_genes</code> output.
<code>explicit.title</code>	logical specifying whether the title should include <code>logFC</code> information for each condition.
<code>RP.MT.filter</code>	Boolean. Defaults to <code>TRUE</code> . Whether to exclude ribosomal and mitochondrial genes.
<code>color.p.threshold</code>	numeric specifying the adjusted p-value threshold for <code>geom_points</code> to be colored. Default is set to 0.01.
<code>color.log.threshold</code>	numeric specifying the absolute <code>logFC</code> threshold for <code>geom_points</code> to be colored. Default is set to 0.25.
<code>label.p.threshold</code>	numeric specifying the adjusted p-value threshold for genes to be labeled via <code>geom_text_repel</code> . Default is set to 0.001.
<code>label.logfc.threshold</code>	numeric specifying the absolute <code>logFC</code> threshold for genes to be labeled via <code>geom_text_repel</code> . Default is set to 0.75.
<code>n.label.up</code>	numeric specifying the number of top upregulated genes to be labeled via <code>geom_text_repel</code> . Genes will be ordered by adjusted p-value. Overrides the " <code>label.p.threshold</code> " and " <code>label.logfc.threshold</code> " parameters.
<code>n.label.down</code>	numeric specifying the number of top downregulated genes to be labeled via <code>geom_text_repel</code> . Genes will be ordered by adjusted p-value. Overrides the " <code>label.p.threshold</code> " and " <code>label.logfc.threshold</code> " parameters.
<code>by.logFC</code>	logical. If set to <code>TRUE</code> <code>n.label.up</code> and <code>n.label.down</code> will label genes ordered by <code>logFC</code> instead of adjusted p-value.
<code>maximum.overlaps</code>	integer specifying removal of labels with too many overlaps. Default is set to <code>Inf</code> .
<code>plot.adj.pvalue</code>	logical specifying whether adjusted p-value should be plotted on the y-axis.

Value

Returns a volcano plot from the output of the `FindMarkers` function from the `Seurat` package, which is a `ggplot` object that can be modified or plotted. Infinite p-values are set defined value of the highest $-\log(p) + 100$.

Examples

```
## Not run:
#using the findmarkers.output
GEX_volcano(findmarkers.output = FindMarkers.Output
```

```
, condition.1 = "cluster1", condition.2 = "cluster2"
, maximum.overlaps = 20)

GEX_volcano(findmarkers.output = FindMarkers.Output
, condition.1 = "cluster1", condition.2 = "cluster2"
, n.label.up = 50, n.label.down = 20)

#using the GEX_cluster_genes output
GEX_volcano(findmarkers.output = GEX_cluster_genes.Output
, cluster.genes.output =TRUE)

## End(Not run)
```

hotspot_df

hotspot_df Hotspot mutations taken from Yaari et al., *Frontiers in Immunology*, 2013. This contains transition probabilities for all 5mer combinations based on high throughput sequencing data. The transition probabilities are for the middle nucleotide in each 5mer set. This can be customized by changing the genes and sequences. Custom mutation hotspots can be supplied by modifying this dataframe. Repeating particular hotspot entries allows for the hotspot to mutate more than one time per SHM event.

Description

@format A data frame with 1024 rows and 6 variables:

pattern Character array where each entry corresponds to a 5 base motif. The mutation probabilities correspond to the middle nucleotide in each 5mer.

toA The probability for the middle nucleotide in "pattern" to mutate to an adenine

toC The probability for the middle nucleotide in "pattern" to mutate to an cytosine

toG The probability for the middle nucleotide in "pattern" to mutate to an guanine

toT The probability for the middle nucleotide in "pattern" to mutate to an thymine

Source The origin of how this motif was discovered. Either Inferred or Experimental

Usage

```
data("hotspot_df")
```

Format

An object of class `data.frame` with 1024 rows and 6 columns.

Source

Yaari et al., *Frontiers in Immunology*, 2013

hum_b_h	<i>hum_b_h</i>
---------	----------------

Description

human germline IgH (heavy chain v,d,j gene segments). When multiple alleles were present, the first one was included. These names and sequences can be changed by customized by changing this dataframe. Additionally, repeating elements can give certain germline gene elements a larger probability of being used during repertoire evolution.

Usage

```
data("hum_b_h")
```

Format

A list including 3 elements (data frames): v gene, d gene, j gene, respectively.

```
[[1]]
```

gene The v gene name

seq The corresponding sequence [[2]]

gene The d gene name

seq The corresponding sequence [[3]]

gene The j gene name

seq The corresponding sequence

Source

IMGT

hum_b_l	<i>hum_b_l</i>
---------	----------------

Description

human germline IgH (light chain v,d,j gene segments). When multiple alleles were present, the first one was included. These names and sequences can be changed by customized by changing this dataframe. Additionally, repeating elements can give certain germline gene elements a larger probability of being used during repertoire evolution.

Usage

```
data("hum_b_l")
```

Format

A list including 2 elements (data frames): v gene, d gene, j gene, respectively.

[[1]]

gene The v gene name

seq The corresponding sequence [[2]]

gene The j gene name

seq The corresponding sequence

Source

IMGT

hum_t_h

hum_t_h

Description

human germline TRB (heavy chain v,d,j gene segments). When multiple alleles were present, the first one was included. These names and sequences can be changed by customized by changing this dataframe. Additionally, repeating elements can give certain germline gene elements a larger probability of being used during repertoire evolution.

Usage

```
data("hum_t_h")
```

Format

A list including 3 elements (data frames): v gene, d gene, j gene, respectively.

[[1]]

gene The v gene name

seq The corresponding sequence [[2]]

gene The d gene name

seq The corresponding sequence [[3]]

gene The j gene name

seq The corresponding sequence

Source

IMGT

hum_t_1	<i>hum_t_1</i>
---------	----------------

Description

human germline TRA (light chain v,d,j gene segments). When multiple alleles were present, the first one was included. These names and sequences can be changed by customized by changing this dataframe. Additionally, repeating elements can give certain germline gene elements a larger probability of being used during repertoire evolution.

Usage

```
data("hum_t_1")
```

Format

A list including 2 elements (data frames): v gene, d gene, j gene, respectively.

```
[[1]]
```

gene The v gene name

seq The corresponding sequence [[2]]

gene The j gene name

seq The corresponding sequence

Source

IMGT

iso_SHM_prob	<i>iso_SHM_prob</i> A probability dataframe specifying SHM.nuc.prob for cells of different isotypes. The first column is the names of isotypes, while the second column is the SHM.nuc.prob of cell of that isotype. user can define different SHM.nuc.prob for isotypes.
--------------	---

Description

iso_SHM_prob A probability dataframe specifying SHM.nuc.prob for cells of different isotypes. The first column is the names of isotypes, while the second column is the SHM.nuc.prob of cell of that isotype. user can define different SHM.nuc.prob for isotypes.

Usage

```
data("iso_SHM_prob")
```

Format

a dataframe with 2 columns

 mus_b_h

mus_b_h

Description

C57BL/6 germline IgH (heavy chain v,d,j gene segments). When multiple alleles were present, the first one was included. These names and sequences can be changed by customized by changing this dataframe. Additionally, repeating elements can give certain germline gene elements a larger probability of being used during repertoire evolution.

Usage

```
data("mus_b_h")
```

Format

A list including 3 elements (data frames): v gene, d gene, j gene, respectively.

```
[[1]]
```

gene The v gene name

seq The corresponding sequence [[2]]

gene The d gene name

seq The corresponding sequence [[3]]

gene The j gene name

seq The corresponding sequence

Source

IMGT

 mus_b_l

mus_b_l

Description

C57BL/6 germline IgH (light chain v,d,j gene segments). When multiple alleles were present, the first one was included. These names and sequences can be changed by customized by changing this dataframe. Additionally, repeating elements can give certain germline gene elements a larger probability of being used during repertoire evolution.

Usage

```
data("mus_b_l")
```

Format

A list including 2 elements (data frames): v gene, d gene, j gene, respectively.

[[1]]

gene The v gene name

seq The corresponding sequence [[2]]

gene The j gene name

seq The corresponding sequence

Source

IMGT

mus_b_trans

mus_b_trans A data frame contains mouse B cell average gene expression for multiple cell types, with the rows representing the gene names, column names representing the cell type names. The original single cell sequencing data is retrieved from 10xgenomics and combined with experimental data from.#? The expression level for different cell types are obtained by calculating the average expression after sorting the original data by markers as shown below.

NaiveBcell Cd19+;Cd27-;Cd38-

GerminalcenterBcell Fas+;Cd19+

Plasmacell Sdc1+

MemoryBcell Cd38+;Fas-

Description

mus_b_trans A data frame contains mouse B cell average gene expression for multiple cell types, with the rows representing the gene names, column names representing the cell type names. The original single cell sequencing data is retrieved from 10xgenomics and combined with experimental data from.#? The expression level for different cell types are obtained by calculating the average expression after sorting the original data by markers as shown below.

NaiveBcell Cd19+;Cd27-;Cd38-

GerminalcenterBcell Fas+;Cd19+

Plasmacell Sdc1+

MemoryBcell Cd38+;Fas-

Usage

data("mus_b_trans")

Format

A data frame with 26538 rows and 4 variables, with the rows representing the gene names, column names representing the cell type names.

Source

https://support.10xgenomics.com/single-cell-vdj/datasets/3.0.0/vdj_v1_mm_c57bl6_pbmc_5gex https://support.10xgenomics.com/single-cell-vdj/datasets/3.0.0/vdj_v1_mm_balbc_pbmc_5gex

mus_t_h

mus_t_h

Description

C57BL/6 germline TRB (heavy chain v,d,j gene segments). When multiple alleles were present, the first one was included. These names and sequences can be changed by customized by changing this dataframe. Additionally, repeating elements can give certain germline gene elements a larger probability of being used during repertoire evolution.

Usage

```
data("mus_t_h")
```

Format

A list including 3 elements (data frames): v gene, d gene, j gene, respectively.

```
[[1]]
```

gene The v gene name

seq The corresponding sequence [[2]]

gene The d gene name

seq The corresponding sequence [[3]]

gene The j gene name

seq The corresponding sequence

Source

IMGT

`mus_t_1`*mus_t_1*

Description

C57BL/6 germline TRA (light chain v,d,j gene segments). When multiple alleles were present, the first one was included. These names and sequences can be changed by customized by changing this dataframe. Additionally, repeating elements can give certain germline gene elements a larger probability of being used during repertoire evolution.

Usage

```
data("mus_t_1")
```

Format

A list including 2 elements (data frames): v gene, d gene, j gene, respectively.

```
[[1]]
```

gene The v gene name

seq The corresponding sequence [[2]]

gene The j gene name

seq The corresponding sequence

Source

IMGT

`new`*Example repertoire for testing and examples*

Description

Example repertoire for testing and examples

Usage

```
new
```

Format

An object of class `list` of length 1662.

no.empty.node	<i>Get clone network igraphs without empty mode. Empty node represents the 'extincted' sequences, that are not in any living cell but once existed.</i>
---------------	---

Description

Get clone network igraphs without empty mode. Empty node represents the 'extincted' sequences, that are not in any living cell but once existed.

Usage

```
no.empty.node(history, igrph.index)
```

Arguments

history	The dataframe 'history' from the simulation output.
igrph.index	The list 'igrph.index' from the simulation output.
empty.node	If TRUE, there will be empty node in igrph. if FALSE, the empty node will be deleted.

Value

a list of clone network igraphs without empty mode.

one_spot_df	<i>one_spot_df</i>
-------------	--------------------

Description

WRC hotspot mutations taken from Yaari et al., *Frontiers in Immunology*, 2013. These include only the mutations following the WRC pattern, where W equals A or T and R equals A or G). Custom mutation hotspots can be supplied by modifying this dataframe. Repeating particular hotspot entries allows for the hotspot to mutate more than one time per SHM event.

Usage

```
data("one_spot_df")
```


Format

A data frame with 32 rows and 6 variables:

pattern Character array where each entry corresponds to a 5 base motif. The mutation probabilities correspond to the middle nucleotide in each 5mer.

toA The probability for the middle nucleotide in "pattern" to mutate to an adenine

toC The probability for the middle nucleotide in "pattern" to mutate to an cytosine

toG The probability for the middle nucleotide in "pattern" to mutate to an guanine

toT The probability for the middle nucleotide in "pattern" to mutate to an thymine

Source The origin of how this motif was discovered. Either Inferred or Experimental

Source

Yaari et al., Frontiers in Immunology, 2013

pheno_SHM_prob	<i>pheno_SHM_prob A probability dataframe specifying SHM.nuc.prob for cells of different phenotypes. The first column is the names of phenotypes, while the second column is the SHM.nuc.prob of cell of that phenotype. user can define different SHM.nuc.prob for phenotypes.</i>
----------------	---

Description

pheno_SHM_prob A probability dataframe specifying SHM.nuc.prob for cells of different phenotypes. The first column is the names of phenotypes, while the second column is the SHM.nuc.prob of cell of that phenotype. user can define different SHM.nuc.prob for phenotypes.

Usage

```
data("pheno_SHM_prob")
```

Format

a dataframe with 2 columns

PlatypusDB_AIRR_to_VGM

Loads in and converts input AIRR-compatible tsv file(s) into the Platypus VGM object format. All compulsory AIRR data columns are needed. Additionally, the following columns are required: c_call, cell_id, clone_id. If trim.and.align is set to TRUE additionally the following columns are needed: v_sequence_start, j_sequence_end.

Description

Loads in and converts input AIRR-compatible tsv file(s) into the Platypus VGM object format. All compulsory AIRR data columns are needed. Additionally, the following columns are required: c_call, cell_id, clone_id. If trim.and.align is set to TRUE additionally the following columns are needed: v_sequence_start, j_sequence_end.

Usage

```
PlatypusDB_AIRR_to_VGM(
  AIRR.input,
  get.VDJ.stats,
  VDJ.combine,
  trim.and.align,
  filter.overlapping.barcodes.VDJ,
  group.id,
  verbose
)
```

Arguments

AIRR.input	Source of the AIRR table(s) as a list. There are 2 available input options: 1. List with local paths to .tsv files / 3. List of AIRR tables loaded in as R objects within the current R environment.
get.VDJ.stats	Boolean. Defaults to TRUE. Whether to generate summary statistics on repertoires and output those as output_VGM[[3]]
VDJ.combine	Boolean. Defaults to TRUE. Whether to integrate repertoires. A sample identifier will be appended to each barcode both. Highly recommended for all later functions
trim.and.align	Boolean. defaults to FALSE. Whether to trim VJ/VDJ seqs and add information from alignment in AIRR dataframe columns. ! No alignment is done here, instead, columns containing alignment information in the AIRR dataframes are reformatted.
filter.overlapping.barcodes.VDJ	Boolean. defaults to TRUE. Whether to remove barcodes which are shared among samples in the GEX analysis. Shared barcodes normally appear at a very low rate.

group.id	vector with integers specifying the group membership. c(1,1,2,2) would specify the first two elements of the input AIRR list are in group 1 and the third/fourth input elements will be in group 2.
verbose	Writes runtime status to console. Defaults to FALSE

Value

A VDJ_GEX_Matrix object used in Platypus V3 as an input to most analysis and plotting functions

Examples

```
## Not run:

VGM <- PlatypusDB_AIRR_to_VGM(AIRR.input =
list("~/pathto/s1/airr_rearrangement.tsv", "~/pathto/s2/airr_rearrangement.tsv"),
VDJ.combine = TRUE, group.id = c(1,2), filter.overlapping.barcodes.VDJ = TRUE)

## End(Not run)
```

PlatypusDB_fetch	<i>Loads and saves RData objects from the PlatypusDB</i>
------------------	--

Description

Loads and saves RData objects from the PlatypusDB

Usage

```
PlatypusDB_fetch(
  PlatypusDB.links,
  save.to.disk,
  load.to.environment,
  load.to.list,
  path.to.save,
  combine.objects
)
```

Arguments

PlatypusDB.links

Character vector. One or more links to files in the PlatypusDB. Links are constructed as follows: "%Project id%/%sample_id%/%filetype%". Any of the three can be "ALL", to download all files fitting the other link elements. If %filetype% is gexVGM, vjVGM or metadata, %sample_id% needs to be "ALL", as these are elements which are not divided by sample. See examples for clarification. See last example on how to download AIRR compliant data. Feature Barcode (FB) data will be downloaded both for GEX and VDJ if present and does not need to be specified in the path

save.to.disk	Boolean. Defaults to FALSE. Whether to save downloaded files individually to the directory specified in path.to.save
load.to.enviroment	Boolean. Defaults to TRUE. Whether to load objects directly into the current .GlobalEnv. An array of the names of the loaded objects will be returned. !Be aware of RAM limitations of your machine when downloading multiple large files.
load.to.list	Boolean. Defaults to FALSE. Whether to return loaded objects as a list. !Be aware of RAM limitations of your machine when downloading multiple large files.
path.to.save	System path to save files to.
combine.objects	Boolean. Defaults to TRUE. Whether to combine objects if appropriate. e.g. VDJ and GEX RData objects for a sample are saved as two independent objects and downloaded as such, to allow for flexibility. If combine.objects is set to TRUE, the function will coerce RData objects of each loaded sample or of each loaded VDJ_GEX_matrix appropriately. Combined input of VDJ and GEX Rdata objects can be directly supplied to the VDJ_GEX_matrix function.

Value

A list of loaded project files as R objects if load.to.list = T or a name of these object loaded to the enviroment if load.to.enviroment = T.

Examples

```
#Get a list of available projects by name
names(PlatypusDB_list_projects())

#Load the VDJ_GEX_matrix of a project as an object and
#also save it to disk for later.
#This will download the VDJ and GEX part of the VDJ_GEX_matrix and combine
PlatypusDB_fetch(PlatypusDB.links = c("Kuhn2021a//ALL")
,save.to.disk = FALSE,load.to.enviroment = TRUE, load.to.list = FALSE
, combine.object = TRUE,path.to.save = "/Downloads")

#Load VDJ dataframe of the VDJ GEX matrix for all samples of one project
loaded_list <- PlatypusDB_fetch(PlatypusDB.links = c("Kuhn2021a//VDJmatrix")
,save.to.disk = FALSE,load.to.enviroment = FALSE, load.to.list = TRUE)

#Load the VDJ and GEX RData of 2 samples from
#2 different projects which can be directly passed
#on to the VDJ_GEX_matrix function to integrate
#downloaded_objects <- PlatypusDB_fetch(
#PlatypusDB.links = c("Project1/s1/ALL", "Project1/s2/ALL")
#,save.to.disk = FALSE,load.to.enviroment = FALSE, load.to.list = TRUE
#, combine.objects = TRUE)
```

```

#integrated_samples <- VDJ_GEX_matrix_DB(data.in = downloaded_objects)

#Download metadata objects for projects
list_of_metadata_tables <- PlatypusDB_fetch(
PlatypusDB.links = c("Kuhn2021a//metadata")
,save.to.disk = FALSE,load.to.environment = FALSE, load.to.list = TRUE)

#Download of airr_rearrangement.tsv
#Load VDJ.RData into a list
#downloaded_objects <- PlatypusDB_fetch(
#PlatypusDB.links = c("Project1/ALL/VDJ.RData"),save.to.disk = FALSE
#,load.to.environment = FALSE, load.to.list = TRUE)

#Extract airr_rearrangement table for sample 1
#airr_rearrangement <- downloaded_objects[[1]][[1]][[6]]
#Index hierarchy: Sample, VDJ or GEX, VDJ element

#Save for import to AIRR compatible pipeline
#write.table(airr_rearrangement, file = "airr_rearrangement_s1.tsv", sep='\t',
#row.names = FALSE, quote=FALSE)

```

PlatypusDB_find_CDR3s *Queries for the occurrence of CDR3 sequences in public datasets on PlatypusDB.*

Description

Queries for the occurrence of CDR3 sequences in public datasets on PlatypusDB.

Usage

```
PlatypusDB_find_CDR3s(VDJ.cdr3s.aa, VJ.cdr3s.aa, projects.to.search)
```

Arguments

VDJ.cdr3s.aa Character A VDJ CDR3s amino acid sequence to search for
VJ.cdr3s.aa Character A VJ CDR3s amino acid sequence to search for
projects.to.search Optional character vector. Defaults to "ALL". Names of projects to search within.

Value

A list of subsets of VDJ matrices from projects containing the query VDJ CDR3 (out[[1]]), the VJ CDR3 (out[[2]]) and cells containing both the query VDJ and VJ CDR3s (out[[3]])

Examples

```
## Not run:
public_clones <- PlatypusDB_find_CDR3s(VDJ.cdr3s.aa = "CMRYGNYWYFDVW"
, VJ.cdr3s.aa = "CLQHGESPTF", projects.to.search = "ALL")

## End(Not run)
```

PlatypusDB_list_projects

Lists metadata tables of available projects on PlatypusDB

Description

Lists metadata tables of available projects on PlatypusDB

Usage

```
PlatypusDB_list_projects(keyword)
```

Arguments

keyword Character. Keyword by which to search project ids (First Author, Year) in the database. Defaults to an empty string ("") which will list all projects currently available

Value

A list of metadata tables by project. List element names correspond to project ids to use in the PlatypusDB_fetch function

Examples

```
#Get list of all available projects and metadata.
PlatypusDB_projects <- PlatypusDB_list_projects()

#Names of list are project ids to use in PlatypusDB_fetch function
names(PlatypusDB_projects)
#Common format: first author, date, letter a-z (all lowercase)

#View metadata of a specific project
print(PlatypusDB_projects[["Kuhn2021a"]])
```

PlatypusDB_load_from_disk

Utility function for loading in local dataset as VDJ_GEX_matrix and PlatypusDB compatible R objects. Especially useful when wanting to integrate local and public datasets. This function only imports and does not make changes to format, row and column names. Exception: filtered_contig.fasta are appended to the filtered_contig_annotations.csv as a column for easy access

Description

Utility function for loading in local dataset as VDJ_GEX_matrix and PlatypusDB compatible R objects. Especially useful when wanting to integrate local and public datasets. This function only imports and does not make changes to format, row and column names. Exception: filtered_contig.fasta are appended to the filtered_contig_annotations.csv as a column for easy access

Usage

```
PlatypusDB_load_from_disk(
  VDJ.out.directory.list,
  GEX.out.directory.list,
  FB.out.directory.list,
  batches
)
```

Arguments

VDJ.out.directory.list

List containing paths to VDJ output directories from cell ranger. This pipeline assumes that the output file names have not been changed from the default 10x settings in the /outs/ folder. This is compatible with B and T cell repertoires (both separately and simultaneously).

GEX.out.directory.list

List containing paths the outs/ directory of each sample or directly the raw or filtered_feature_bc_matrix folder. Order of list items must be the same as for VDJ. This outs directory may also contain Feature Barcode (FB) information. Do not specify FB.out.directory in this case.

FB.out.directory.list

List of paths pointing at the outs/ directory of output of the Cellranger counts function which contain Feature barcode counts. Any input will overwrite potential FB data loaded from the GEX input directories. Length must match VDJ and GEX directory inputs. (in case of a single FB output directory for multiple samples, please specify this directory as many times as needed)

batches

Integer vector. Defaults to all 1, yielding all samples with batch number "b1". Give a batch number to each sample (each entry in the VDJ/GEX input lists). This will be saved as element 5 in the sample list output.

Value

Large nested list object containing all needed Cellranger outputs to run the VDJ_GEX_matrix function. Level 1 of the list are samples, level 2 are VDJ GEX and metadata information. (e.g. out[[1]][[1]] corresponds to VDJ data objects of sample 1)

Examples

```
## Not run:
VDJ.in <- list()
VDJ.in[[1]] <- c("~/VDJ/S1/")
VDJ.in[[2]] <- c("~/VDJ/S2/")
GEX.in <- list()
GEX.in[[1]] <- c("~/GEX/S1/")
GEX.in[[2]] <- c("~/GEX/S2/")
PlatypusDB_load_from_disk(VDJ.out.directory.list = VDJ.in, GEX.out.directory.list = GEX.in)

## End(Not run)
```

PlatypusDB_VGM_to_AIRR

Exports AIRR compatible tables supplemented with VDJ and GEX information from the Platypus VGM object and the cellranger output airr_rearrangements.tsv

Description

Exports AIRR compatible tables supplemented with VDJ and GEX information from the Platypus VGM object and the cellranger output airr_rearrangements.tsv

Usage

```
PlatypusDB_VGM_to_AIRR(
  VGM,
  VDJ.features.to.append,
  GEX.features.to.append,
  airr.rearrangements
)
```

Arguments

VGM Output object of the VDJ_GEX_matrix function generated with VDJ.combine = T, GEX.combine = T (to merge all samples) and integrate.VDJ.to.GEX = T (to integrate VDJ and GEX data)

VDJ.features.to.append Column names of the VGM VDJ matrix (VGM[[1]]) to append to the AIRR compatible table

GEX.features.to.append

Metadata column names or Gene names of the VGM GEX object (VGM[[2]]) to append to the AIRR compatible table. For a list of available features run: names(VGM[[2]]@meta.data) and rownames(VGM[[2]])

airr.rearrangements

Source of the airr_rearrangements.tsv file as generated by cellranger. There are 3 available input options: 1. R list object from Platypus_DB_load_from_disk or Platypus_DB_fetch / 2. List with local paths to airr_rearrangements.tsv / 3. List of airr_rearrangements.tsv loaded in as R objects within the current R environment. ! Order of input list must be identical to that of sample_ids in the VGM !

Value

A list of length of samples in VGM containing a AIRR-compatible dataframe for each sample. ! Cave the format: VGM object => 1 cell = 1 row; AIRR table 1 cell = as many rows as VDJ and VJ chains available for that cell. GEX cell-level information is attached to all rows containing a chain of that cell.

Examples

```
## Not run:
#complete workflow below
#usage with airr rearrangement tables from PlatypusDB_load_from_disk
#or PlatypusDB_fetch list object
airr.list.out <- PlatypusDB_VGM_to_AIRR(VGM = VGM
, VDJ.features.to.append = c("VDJ_cdr3s_aa")
, GEX.features.to.append = c("CTLA4", "TOX"), airr.rearrangements = Data.in)

#usage with airr rearrangement tables from disk
airr.list.out <- PlatypusDB_VGM_to_AIRR(VGM = VGM
, VDJ.features.to.append = c("VDJ_cdr3s_aa")
, GEX.features.to.append = c("CTLA4", "TOX"),
airr.rearrangements =list("~/path_to/s1/airr_rearrangement.tsv"
, "~/path_to/s2/airr_rearrangement.tsv"))

#usage with airr rearrangement tables from objects in R environment
airr.list.out <- PlatypusDB_VGM_to_AIRR(VGM = VGM
, VDJ.features.to.append = c("VDJ_cdr3s_aa")
, GEX.features.to.append = c("CTLA4", "TOX"),
airr.rearrangements = list(airr_rearrangements.s1, airr_rearrangements.s2))

#Complete workflow
#set paths of cellranger directories containing
#also the airr_rearrangements.tsv file
VDJ.out.directory.list <- list()
VDJ.out.directory.list[[1]] <- c("~/cellrangerVDJ/s1")
VDJ.out.directory.list[[2]] <- c("~/cellrangerVDJ/s2")

GEX.out.directory.list <- list()
GEX.out.directory.list[[1]] <- c("~/cellrangerGEX/s1")
GEX.out.directory.list[[2]] <- c("~/cellrangerGEX/s2")
```

```

#Run VGM with GEX and VDJ integration
VGM <- VDJ_GEX_matrix(VDJ.out.directory.list = VDJ.out.directory.list,
GEX.out.directory.list = GEX.out.directory.list,
GEX.integrate = T, VDJ.combine = T, integrate.GEX.to.VDJ = T
, integrate.VDJ.to.GEX = T,
get.VDJ.stats = F, trim.and.align = F)
#Generate AIRR compatible table supplemented by GEX information
airr.list.out <- PlatypusDB_VGM_to_AIRR(VGM = VGM,
VDJ.features.to.append = c("VDJ_sequence_nt_trimmed", "VJ_sequence_nt_trimmed"),
GEX.features.to.append = c("UMAP_1", "UMAP_2", "CTLA4", "TOX"),
airr.rearrangements = c("~/cellrangerVDJ/s1/airr_rearrangement.tsv"
, "~/cellrangerVDJ/s2/airr_rearrangement.tsv"))

#To save a dataframe as .tsv
write.table(airr_dataframe, file = "supplemented_airr_rearrangements.tsv"
, sep='\t', row.names = FALSE, quote=FALSE)

## End(Not run)

```

select.top

Get the index of top ranking clones.

Description

Get the index of top ranking clones.

Usage

```
select.top(clonotypes, top.n)
```

Arguments

clonotypes	The output "clonotypes" dataframe from simulation output.
top.n	The top n abundant clones to be selected.

Value

a vector of indexes of top ranking clones

simulate_repertoire *Simulate immune repertoire and transcriptome data*

Description

Simulate repertoire and transcriptome matrix, with igraph tree plot for each clone showing the evolution process. the node in the tree plot are colored with transcriptome state and isotype.

Usage

```
simulate_repertoire(  
  initial.size.of.repertoire,  
  species,  
  cell.type,  
  cd4.proportion,  
  duration.of.evolution,  
  complete.duration,  
  vdj.productive,  
  vdj.model,  
  vdj.insertion.mean,  
  vdj.insertion.stdv,  
  vdj.branch.prob,  
  clonal.selection,  
  cell.division.prob,  
  sequence.selection.prob,  
  special.v.gene,  
  class.switch.prob,  
  class.switch.selection.dependent,  
  class.switch.independent,  
  SHM.method,  
  SHM.nuc.prob,  
  SHM.isotype.dependent,  
  SHM.phenotype.dependent,  
  max.cell.number,  
  max.clonotype.number,  
  death.rate,  
  igraph.on,  
  transcriptome.on,  
  transcriptome.switch.independent,  
  transcriptome.switch.prob,  
  transcriptome.switch.isotype.dependent,  
  transcriptome.switch.SHM.dependent,  
  transcriptome.switch.selection.dependent,  
  transcriptome.states,  
  transcriptome.noise,  
  seq.name  
)
```

Arguments

<code>initial.size.of.repertoire</code>	The initial number of existing cells when the evolution starts. Default is 10.
<code>species</code>	The species of the simulated repertoire, can be "mus" for mouse or "hum" for human. Default is "mus".
<code>cell.type</code>	The cell type for the simulation. "B" or "T"
<code>cd4.proportion</code>	A number between 0 and 1 specifying the proportion of Cd4+ T cells, when <code>cell.type</code> is "T" and <code>transcripttime</code> states data is default. Default is 1, all the cells are Cd4. When user specify transcriptome data for T cells, mixture of CD4+ and CD8+ T cells are not applicable.
<code>duration.of.evolution</code>	The maxim time steps for simulation.
<code>complete.duration</code>	TRUE or FALSE. Default is TURE. If TURE, after cell number or clone number reaches the upper limit, the evolution(class switch, mutation, transcriptional state switch) will continue until the <code>duration.of.evolution</code> is complete. If FLASE, the evolution will stop when either cell number or clone number reaches the limit.
<code>vdj.productive</code>	"random": the sequence will be generated from random VDJ recombination, there might be a proportion of unproductive sequences. These VDJ genes were taken from IMGT. When more than one allele was present for a given gene, the first was used. "naive": the VDJ sequence be sampled from a pool of productive sequences obtained by filtering randomly simulated sequences with MIXCR. "vae": the VDJ sequence be sampled from a pool of productive sequences obtained by filtering sequences generated from vae models with MIXCR.
<code>vdj.model</code>	Specifies the model used to simulate V-D-J recombination. Can be either "naive" or "data". "naive" is chain independent and does not differentiate between different species. To rely on the default "experimental" options, this should be "data" and the parameter <code>vdj.insertion.mean</code> should be "default". This will allow for different mean additions for either the VD and JD junctions and will differ depending on species.
<code>vdj.insertion.mean</code>	Integer value describing the mean number of nucleotides to be inserted during simulated V-D-J recombination events. If "default" is entered, the mean will be normally distributed.
<code>vdj.insertion.stdv</code>	Integer value describing the standard deviation corresponding to insertions of V-D-J recombination. No "default" parameter currently supported but will be updated with future experimental data. This should be a number if using a custom distribution for V-D-J recombination events, but can be "default" if using the "naive" <code>vdj.model</code> or the "data", with <code>vdj.insertion.mean</code> set to "default".
<code>vdj.branch.prob</code>	Probability of new VDJ recombination event in each time step. when new VDJ recombination happen, a new cell with a new sequence will be generated. Default is 0.2.

<code>clonal.selection</code>	TRUE or FALSE. If TRUE, cells in clones with higher frequency have their division probability proportional to the clonal frequency. If FALSE, clones with higher frequency will have lower probability to expand.
<code>cell.division.prob</code>	Probability of cells to be duplicated in each time step. Default is 0.1. If uneven probability for different clones is needed, the input should be a vector of 2 numeric items, with the first item being the lower bound, the second item being the upper bound of the division rate. The most abundant clone will get the highest division rate, and division rate of other clones will follow arithmetic progression and keep decreasing until the last abundant clone with the lower limit of division rate. If input 3 values, the third value will be the division rate for cells with selected sequences. If a fourth number is given, the division probability of selected sequence will be sampled between the third number and the fourth number.
<code>sequence.selection.prob</code>	Probability of each unique sequence to be selected as expanding sequence. Expanding sequences can have their division rate specified in the third element of <code>cell.division.prob</code> .
<code>special.v.gene</code>	If TRUE, simulation will apply <code>sequence.selection.prob</code> for heavy and light chain v gene combination specified in dataframe "special_v".
<code>class.switch.prob</code>	Probability matrix of class switching for b cells. The row names of the matrix are the isotypes the cell is switching from, the column names are the isotypes the cell is switching to. All B cells start from IGHM, and switch to one of the other isotypes or remain the same. Default values are in the attaching matrix "class_switch_prob_hum" and "class_switch_prob_mus". The order of isotype in rows and columns should be the same.
<code>class.switch.selection.dependent</code>	If TRUE, class switching will happen when the cell is selected, if the cell has IgM or IgD isotype.
<code>class.switch.independent</code>	If TRUE, class switching will happen randomly at each time step for all cells. If FALSE, random class switching will be switched off.
<code>SHM.method</code>	The mode of SHM speciation events. Options are either: "poisson", "data", "motif", "wrc", and "all". Specifying either "poisson" or "naive" will result in mutations that can occur anywhere in the heavy chain region, with each nucleotide having an equal probability for a mutation event. Specifying "data" focuses mutation events during SHM in the CDR regions (based on IMGT), and there will be an increased probability for transitions (and decreased probability for transversions). Specifying "motif" will cause neighbor dependent mutations based on a mutational matrix from high throughput sequencing data sets (Yaari et al., <i>Frontiers in Immunology</i> , 2013). "wrc" allows for only the WRC mutational hotspots to be included (where W equals A or T and R equals A or G). Specifying "all" will use all four types of mutations during SHM branching events, where the weights for each can be specified in the "SHM.nuc.prob" parameter.
<code>SHM.nuc.prob</code>	Specifies the rate at which nucleotides change during speciation (SHM) events. This parameter depends on the type of mutation specified by <code>SHM.method</code> . For

both "poisson" and "data", the input value determines the probability for each site to mutate (the whole sequence for "poisson" and the CDRs for "data"). For either "motif" or "wrc", the number of mutations per speciation event should be specified. Note that these are not probabilities, but the number of mutations that can occur (if the mutation is present in the sequence). If "all" is specified, the input should be a vector where the first element controls the poisson style mutations, second controls the "data", third controls the "motif" and fourth controls the "wrc".

- `SHM.isotype.dependent`
If TRUE, somatic hypermutation of certain isotype will happen based on probability specified in dataframe "iso_SHM_prob".
- `SHM.phenotype.dependent`
If TRUE, somatic hypermutation of certain phenotype will happen based on probability specified in dataframe "pheno_SHM_prob".
- `max.cell.number`
Integer value describing maximum number of cells allowed. Default is 1500.
- `max.clonotype.number`
Integer value describing maximum number of clones allowed. cell derived from the same mother cell belong to same clone.
- `death.rate`
Probability of cell death happen to each cell in each time step.
- `igraph.on`
If TRUE, mutational network for every B cell clone will be in the output. If False, the igraphs will not be included.
- `transcriptome.on`
If TRUE, the simulation will include transcriptome data. If FALSE, only v_{dj} sequence will be simulated.
- `transcriptome.switch.independent`
TRUE or FALSE value describing whether transcriptome state is allowed to switch independently, not dependent on class switching or somatic hypermutation. If TRUE, `transcriptome.switch.prob` should be specified to control the probability of transcriptome state switching.
- `transcriptome.switch.prob`
Probability of transcriptome state switching independently. Default values are in the attaching matrix "trans_switch_prob_b" and "trans_switch_prob_t". The order of cell type in rows and columns should be the same, and the order of the cell type in the matrix should match cell type names in `transcriptome.states`.
- `transcriptome.switch.isotype.dependent`
TRUE or FALSE value describing whether transcriptome state of a cell is allowed to switch depending on isotype switching. If TRUE, transcriptome state will switch once class switching happens.
- `transcriptome.switch.SHM.dependent`
TRUE or FALSE value describing whether transcriptome state of a cell is allowed to switch depending on somatic hypermutation. If TRUE, transcriptome state will switch once somatic hypermutation happens.
- `transcriptome.switch.selection.dependent`
If TRUE, selected cells will undergo transcriptome state switching if their transcriptome state is 1.

transcriptome.states	A data.frame specifying base gene expression for different cell type, with gene names as row names, cell type names as column names. When missing, a default data.frame will be used. Default data.frame includes "Germinalcenter-Bcell", "NaiveBcell", "Plasmacell", "MemoryBcell" for B cells, and "Naïve Cd4", "ActivatedCd4", "MemoryCd4", "NaiveCd8", "EffectorCd8", "MemoryCd8", "ExhaustedCd8" for T cells. The order of the cell type names in transcriptome.states should match cell type names in the transcriptome.switch.prob matrix.
transcriptome.noise	A character expression specifying the distribution of noise ratio to be multiplied with the base gene expression for each cell. It should be a text expression that generates a numeric vector, which is of the same length as gene names in the transcriptome.state input. Default value is "rnorm(nrow(transcriptome.states), mean = 1, sd = 0.3)".
seq.name	Integer specifies how many top-ranking clones are included in Seq_Name dataframe in the output list for phylogenetic tree plotting in other pipeline. If missing, Seq_Name won't be included in the output.

Value

A list containing the VDJ sequence and corresponding transcriptome data: "all_contig_annotations", "clonotypes", "all_contig", "consensus", "reference", "reference_real", "transcriptome", "igraph_list_iso", "igraph_list_trans",

small_vgm

Small VDJ GEX matrix (VGM) for function testing purposes

Description

Small VDJ GEX matrix (VGM) for function testing purposes

Usage

```
small_vgm
```

Format

An object of class list of length 5.

References

R package Platypus : <https://doi.org/10.1093/nargab/lqab023>

special_v	<i>special_v a dataframe, of heavy and light chain v gene combination and their probability to be selected for expansion.</i>
-----------	---

Description

special_v a dataframe, of heavy and light chain v gene combination and their probability to be selected for expansion.

Usage

```
data("special_v")
```

Format

An object of class data.frame with 5 rows and 3 columns.

trans_switch_prob_b	<i>trans_switch_prob_b The probability for B cell transcriptome states switching. The row names of the matrix are the cell states the cell is switching from, the column names are the cells states the cell is switching to.</i>
---------------------	---

Description

trans_switch_prob_b The probability for B cell transcriptome states switching. The row names of the matrix are the cell states the cell is switching from, the column names are the cells states the cell is switching to.

Usage

```
data("trans_switch_prob_b")
```

Format

A 4*4 matrix. The row and column names are: "GerminalcenterBcell", "NaiveBcell", "Plasmacell", "MemoryBcell". The probability for a cell to switch from "GerminalcenterBcell" to "Plasmacell" is the value at trans_switch_prob_b[1,3].

trans_switch_prob_t *trans_switch_prob_t* The probability for T cell transcriptome states switching. The row names of the matrix are the cell states the cell is switching from, the column names are the cells states the cell is switching to.

Description

trans_switch_prob_t The probability for T cell transcriptome states switching. The row names of the matrix are the cell states the cell is switching from, the column names are the cells states the cell is switching to.

Usage

```
data("trans_switch_prob_t")
```

Format

A 7*7 matrix. The row and column names are: "NaiveCd4", "ActivatedCd4", "MemoryCd4", "NaiveCd8", "EffectorCd8", "Mem

umap.top.highlight *Set idents for top abundant clones in Seurat object, get ready for highlight the top abundant clones in UMAP.*

Description

Set idents for top abundant clones in Seurat object, get ready for highlight the top abundant clones in UMAP.

Usage

```
umap.top.highlight(gex, all.contig.annotations, top.n)
```

Arguments

gex output from get.umap function.
all.contig.annotations The output dataframe all_contig_annotations from simulation.
top.n The top n abundant clones to be shown in the plot. If missing, all clones will be shown.

Value

a Seurat object ready for highlight the top abundant clones in UMAP

VDJ_alpha_beta_Vgene_circos

Produces a Circos plot from the VDJ_analyze output. Connects the V-alpha with the corresponding V-beta gene for each clonotype.

Description

Produces a Circos plot from the VDJ_analyze output. Connects the V-alpha with the corresponding V-beta gene for each clonotype.

Usage

```
VDJ_alpha_beta_Vgene_circos(  
  VDJ,  
  V.or.J,  
  B.or.Tcells,  
  label.threshold,  
  c.threshold,  
  cell.level,  
  clonotype.per.gene.threshold,  
  c.count,  
  platypus.version,  
  filter1H1L  
)
```

Arguments

VDJ	The output of the VDJ_GEX_integrate function (Platypus platypus.version v2). A list of data frames for each sample containing the clonotype information and cluster membership information. For Platypus platypus.version v3, VDJ_GEX_matrix.output[[1]] has to be supplied.
V.or.J	Determines whether to plot the alpha beta gene pairing of the V or J genes. "V", "J" or "both" as possible inputs. Default: "both".
B.or.Tcells	Specify whether B or T cells are being analyzed ("B" or "T"). If not specified, function attempts to decide based on gene names.
label.threshold	Minimal amount of clonotypes per gene necessary to add a gene label to the sector. Default: 0.
c.threshold	Only clonotypes are considered with a frequency higher then c.threshold. Allows to filter for only highly expanded clonotypes.
cell.level	Logical, defines whether weight of connection should be based on number of clonotypes or number of cells. Default: number of clonotypes.
clonotype.per.gene.threshold	How many clonotypes are required to plot a sector for a gene. Filters the rows and columns of the final adjacency matrix.

`c.count` Show clonotype or cell count on Circos plot. Default = T.
`platypus.version` Which platypus.version of platypus is being used. Default = "v3".
`filter1H1L` Whether to filter the input VDJ.matrix in "v3" to only include cells with 1 VDJ and 1 VJ chain. Defaults to TRUE

Value

Returns list of plots. The first n elements contain the circos plot of the n datasets from the VDJ.analyze function. The n+1 element contains a list of the n adjacency matrices for each dataset.

Examples

```
plots <- VDJ_alpha_beta_Vgene_circos(Platypus::small_vgm[[1]]
, platypus.version="v3")
```

VDJ_analyze

Platypus V2 Processes and organizes the repertoire sequencing data from cellranger vdj and returns a list of dataframes, where each dataframe corresponds to an individual repertoire. The function will return split CDR3 sequences, germline gene information, filter out those clones with either incomplete information or doublets (multiple CDR3 sequences for a given chain). This function should be called once for desired integrated repertoire and transcriptome. For example, if there are 3 VDJ libraries and 3 GEX libraries and the goal is to analyze all three GEX libraries together (e.g. one UMAP/tSNE reduction) this then function should be called one time and the three VDJ directories should be provided as input to the single function call.

Description

Platypus V2 Processes and organizes the repertoire sequencing data from cellranger vdj and returns a list of dataframes, where each dataframe corresponds to an individual repertoire. The function will return split CDR3 sequences, germline gene information, filter out those clones with either incomplete information or doublets (multiple CDR3 sequences for a given chain). This function should be called once for desired integrated repertoire and transcriptome. For example, if there are 3 VDJ libraries and 3 GEX libraries and the goal is to analyze all three GEX libraries together (e.g. one UMAP/tSNE reduction) this then function should be called one time and the three VDJ directories should be provided as input to the single function call.

Usage

```
VDJ_analyze(
  VDJ.out.directory,
  filter.1HC.1LC,
  clonotype.list,
```

```

    contig.list,
    filtered.contigs
  )

```

Arguments

VDJ.out.directory
Character vector with each element containing the path to the output of cellranger vdj runs. Multiple repertoires to be integrated in a single transcriptome should be supplied as multiple elements of the character vector. This can be left blank if supplying the clonotypes and contig files directly as input. This pipeline assumes that the output file names have not been changed from the default 10x settings in the /outs/ folder. This is compatible with B and T cell repertoires (both separately and simultaneously).

filter.1HC.1LC Logical indicating whether only those clones containing 1 VH/TRB and VL/TRA should be maintained for further analysis. Default is set to TRUE, which restricts the analysis to only clones with exactly 1 heavy chain and 1 light chain (or 1 beta + 1 alpha in the case of T cells).

clonotype.list List of dataframes containing clonotyping information for each repertoire. The column names should correspond to the clonotypes.csv file from cellranger vdj output.

contig.list List of dataframes containing the contig information for each repertoire. The column names should correspond to the all_contigs.csv file from cellranger vdj output.

filtered.contigs
Logical indicating if the filtered contigs file should be used. TRUE will read VDJ information from only the filtered output of cellranger. FALSE will read the all contigs file from cellranger. Default set to TRUE (filtered output)

Value

Returns a list of dataframes where each dataframe corresponds to one input directory. If only one file is supplied, the output list will only contain one element. This output can be supplied as input to other functions including VDJ_per_clone, VDJ_network, VDJ_germline_genes, VDJ_expansion, visualize_clones_GEX, VDJ_phylo, VDJ_clonotype. Germline gene information is based on the majority of cells within each clonotype. For example, if the majority of cells in clonotype1 have the IGHG1 isotype then the entire clonal family will be determined as IGHG1. For a cell-specific investigation, the output of this function can be supplied to the function VDJ_per_clone, which will provide isotype, sequence, germline gene, etc information for each cell within the each clone.

Examples

```

## Not run:
example.vdj.analyze <- VDJ_analyze(
  VDJ.out.directory = "~/path/to/cellranger/vdj/outs/", filter.1HC.1LC = T)

## End(Not run)

```

VDJ_assemble_for_PnP *Assembles sequences from MIXCR output into inserts for expression in PnP cells. ! ALWAYS VALIDATE INDIVIDUAL SEQUENCE IN GENEIOUS OR OTHER SOFTWARE BEFORE ORDERING SEQUENCES FOR EXPRESSION ! Check notes on column content below ! Only cells with 1 VDJ and 1 VJ sequence are considered. Warnings are issued if sequences do not pass necessary checks*

Description

Assembles sequences from MIXCR output into inserts for expression in PnP cells. ! ALWAYS VALIDATE INDIVIDUAL SEQUENCE IN GENEIOUS OR OTHER SOFTWARE BEFORE ORDERING SEQUENCES FOR EXPRESSION ! Check notes on column content below ! Only cells with 1 VDJ and 1 VJ sequence are considered. Warnings are issued if sequences do not pass necessary checks

Usage

```
VDJ_assemble_for_PnP(  
  VDJ.mixcr.matrix,  
  id.column,  
  species,  
  manual_IgKC,  
  manual_2A,  
  manual_VDJLeader,  
  write.to.disk,  
  filename,  
  verbose  
)
```

Arguments

VDJ.mixcr.matrix	Output dataframe from the VDJ_call_MIXCR function or a dataframe generated using the VDJ_GEX_matrix function and supplemented with MIXCR information (Needed columns: All Framework and CDR sequences)
id.column	Character. Column name of VDJ.mixcr.matrix to use as ID for the assembled sequences. Defaults to "barcode"
species	Character. Which IgKC sequence to use. Can be "human" or "mouse". Defaults to "mouse"
manual_IgKC	Character. Manual overwrite for sequence used as IgKC.
manual_2A	Character. Manual overwrite for sequence used as Furine 2A site.
manual_VDJLeader	Character. Manual overwrite for sequence used as VDJ Leader and signal peptide.

write.to.disk	Boolean. Defaults to TRUE. Whether to save assembled sequences to working directory
filename	Character. Output file name for .fasta and .csv files if write.to.disk == T. Defaults to PnP_assembled_seqs.fasta/csv
verbose	Print runtime message to console. Defaults to FALSE

Value

Returns the input VGM matrix with one additional column containing the assembled sequences. If write.to.disk == T writes a CSV containing key columns of the VGM as well as a .FASTA file to the current working director (getwd()) ! Important notes on column content: 1. The column "seq_length_check" contains either "passed" or "FAILED". If FAILED, this means that at least one of the sequences (e.g. FRL1) was shorter than 9NTs and therefore considered invalid. Please check for missing sequences if you find any warnings 2. The column "seq_codon_check" is deemed "passed" if all CDR and FR input sequences of a cell contain only full codons (i.e. are divisible by 3) 3. The column "PnP_assembled_seqs" contains the assembled sequences / inserts for PnP expression. These should be validated manually in Geneious or other software and can then be ordered to be synthesized. 4. The column "PnP_assembled_annotations" contains a string of annotations for the respective assembled sequence. The structure is | [Sequence element] -> [index (starting from 1) of last nucleotide of the sequence element] ... 5. The column "PnP_assembled_translations" contains the amino acid translation of the full contig that will result from the assembled insert in the backbone PnP vector. Please note: the sequences in the PnP_assembled_translation resulted from pasting the VJ leader sequence (contained in the PnP vector backbone), the PnP_assembled_seqs (The insert itself) and a surrogate stop codon ATAA. If correct, the translation should only contain one * (stop codon) at the very end. For reference: VJLeader sequence: ATGGATTTTCAGGTGCAGATTTTCAGCTTCCTGCTAATCAGCGCTTCAGTTATAATGTCCCGGGGG 6. The column "seq_VJCDR3_check" is deemed "passed" if the translated sequence of the input VJ CDR3 is found in the translated assembled sequence. If this test fails, there is likely an issue with the VJ segment 7. The column "seq_Fur2A_check" is deemed "passed" if correct AA sequence of the 2A site is found in the translated assembled sequence. If this test fails, and the seq_VJCDR3_test was passed, there is likely an issue at the border between VJ and IgKC/2A sequences 8. The column "seq_VDJCDR3_check" is deemed "passed" if the translated sequence of the input VDJ CDR3 is found in the translated assembled sequence. 9. The column "seq_splicesite_check" is deemed passed if the last 6 nucleotides of the assembled sequence are one of the following: "TCCTCA", "TCTTCA", "TCGTCA", "TCATCA".

Examples

```
## Not run:

VGM_with_PnP_seq <- VDJ_assemble_for_PnP(VDJ.mixcr.matrix = VDJ_call_MIXCR.output
, id.column = "barcode", species = "mouse", manual_IgKC = "none", manual_2A = "none"
, manual_VDJLeader = "none", write.to.disk = TRUE, filename = "PnP_seq_example")

## End(Not run)
```

VDJ_call_MIXCR	<i>Extracts information on the VDJRegion level using MiXCR on WINDOWS, MAC and UNIX systems for input from both Platypus v2 (VDJ.per.clone) or v3 (Output of VDJ_GEX_matrix) This function assumes the user can run an executable instance of MiXCR and is eligible to use MiXCR as determined by license agreements. ! FOR WINDOWS USERS THE EXECUTABLE MIXCR.JAR HAS TO PRESENT IN THE CURRENT WORKING DIRECTORY ! The VDJRegion corresponds to the recombined heavy and light chain loci starting from framework region 1 (FR1) and extending to frame work region 4 (FR4). This can be useful for extracting full-length sequences ready to clone and further calculating somatic hypermutation occurrences.</i>
----------------	--

Description

Extracts information on the VDJRegion level using MiXCR on WINDOWS, MAC and UNIX systems for input from both Platypus v2 (VDJ.per.clone) or v3 (Output of VDJ_GEX_matrix) This function assumes the user can run an executable instance of MiXCR and is eligible to use MiXCR as determined by license agreements. ! FOR WINDOWS USERS THE EXECUTABLE MIXCR.JAR HAS TO PRESENT IN THE CURRENT WORKING DIRECTORY ! The VDJRegion corresponds to the recombined heavy and light chain loci starting from framework region 1 (FR1) and extending to frame work region 4 (FR4). This can be useful for extracting full-length sequences ready to clone and further calculating somatic hypermutation occurrences.

Usage

```
VDJ_call_MIXCR(
  VDJ,
  operating.system,
  mixcr.directory,
  species,
  simplify,
  platypus.version
)
```

Arguments

VDJ	For platypus.version = "v2" the output from the VDJ_per_clone function. This object should have information regarding the contigs and clonotype_ids for each cell. For platypus.version = "v3" the VDJ dataframe output of the VDJ_GEX_matrix function (VDJ.GEX.matri.output[[1]])
operating.system	Can be either "Windows", "Darwin" (for MAC) or "Linux". If left empty this is detected automatically
mixcr.directory	The directory containing an executable version of MiXCR. FOR WINDOWS USERS THIS IS SET TO THE CURRENT WORKING DIRECTORY (please

	paste the content of the MIXCR folder after unzipping into your working directory. Make sure, that mixcr.jar is not within any subfolders.)
species	Either "mmu" for mouse or "hsa" for human. These use the default germline genes for both species contained in MIXCR. Defaults to "hsa"
simplify	Only relevant when platypus.version = "v3". Boolean. Defaults to TRUE. If FALSE the full MIXCR output and computed SHM column is appended to the VDJ. If TRUE only the framework and CDR3 region columns and computed SHM column is appended. To discriminate between VDJ and VJ chains, prefixes are added to all MIXCR output columns
platypus.version	Character. Defaults to "v3". Can be "v2" or "v3" dependent on the input format

Value

For platypus.version = "v3" returns input VDJ dataframe supplemented with MIXCR output information. For platypus.version = "v2" returns a nested list containing VDJRegion information as determined by MIXCR. The outer list corresponds to the individual repertoires in the same structure as the input VDJ.per.clone. The inner list corresponds to each clonal family, as determined by either the VDJ_clonotype function or the default nucleotide clonotyping produced by cellranger. Each element in the inner list corresponds to a dataframe containing repertoire information such as isotype, CDR sequences, mean number of UMIs. This output can be supplied to further package functions such as VDJ_extract_sequences and VDJ_GEX_integrate.

See Also

VDJ_extract_sequences

Examples

```
## Not run:
#For platypus version 2
VDJ_call_MIXCR(VDJ = VDJ.per.clone.output,
mixcr.directory = "~/Downloads/mixcr-3.0.12/mixcr", species = "mmu")

#For platypus version 3 on a Windows system
VDJ_call_MIXCR(VDJ = VDJ_GEX_matrix.output[[1]],
mixcr.directory = "WILL BE SET TO CURRENT WORKING DIRECTORY",
species = "mmu", platypus.version = "v3", simplify = TRUE)

## End(Not run)
```

VDJ_circos

Plots a Circos diagram from an adjacency matrix. Uses the Circlize chordDiagram function. Is called by VDJ_clonotype_clusters_circos(), VDJ_alpha_beta_Vgene_circos() and VDJ_VJ_usage_circos() functions or works on its own when supplied with an adjacency matrix.

Description

Plots a Circos diagram from an adjacency matrix. Uses the Circlize chordDiagram function. Is called by VDJ_clonotype_clusters_circos(), VDJ_alpha_beta_Vgene_circos() and VDJ_VJ_usage_circos() functions or works on its own when supplied with an adjacency matrix.

Usage

```
VDJ_circos(Adj_matrix, group, grid.col, label.threshold, axis, c.count)
```

Arguments

Adj_matrix	Adjacency matrix to be plotted. Rownames and Colnames correspond to genes to be matched and entries determine the weight of the connection between the genes (eg. number of clonotypes expressing these two genes).
group	Named list of genes, with list elements corresponding to group-names, and element names being the gene-names. Is generated by VDJ_VJ_usage and VDJ_alpha_beta_Vgene_circos.
grid.col	Named list of genes, with list elements corresponding to color and element names being gene-names. If not supplied it is generated randomly within the function. Is also generated by VDJ_VJ_usage and VDJ_alpha_beta_Vgene_circos.
label.threshold	Genes are only labeled if the count is larger then the label.threshold. By default all label.threshold = 0 (all genes are labeled).
axis	Option to choose the count axis for each gene. "default", "percent" or "max" possible. Default: "max".
c.count	Show clonotype or cell count on Circos plot.

Value

Returns the Circos plot from input of other functions. Do not run as standalone

Examples

```
## Not run:
VDJ_circos() #Do not run as standalone. Called by other circos functions

## End(Not run)
```

VDJ_clonal_donut	<i>Generate circular plots of clonal expansion per repertoire directly from the VDJ matrix of the VDJ_GEX_matrix function</i>
------------------	---

Description

Generate circular plots of clonal expansion per repertoire directly from the VDJ matrix of the VDJ_GEX_matrix function

Usage

```
VDJ_clonal_donut(
  VDJ,
  counts.to.use,
  label.size,
  not.expanded.label.vjust,
  not.expanded.label.hjust,
  total.label.vjust,
  total.label.hjust,
  expanded.colors,
  non.expanded.color
)
```

Arguments

VDJ	VDJ dataframe generated using the VDJ_GEX_matrix function (VDJ_GEX_matrix.output[[1]]). Plots will be made by sample and using the clonal frequencies specified by counts.to.use
counts.to.use	How to count clonotypes and cells. If set to "freq_column" the function uses the clonotype_frequency column derived directly from the cellranger output or the VDJ_clonotype output to calculate expansion. If set to "VGM" the function will base its counts on the number of rows per clonotype in the VDJ. These two counts may diverge, if cells are filtered out due to overlapping barcodes or if a different clonotyping strategy was applied. Defaults to "VGM"
label.size	Size of text labels. All parameters below are purely for graphical purposes and optional. If necessary changes should be made in small (0.1) increments. ! It is recommended to optimize these ONLY once a format for saving the plot is set.
not.expanded.label.vjust	Numeric. Regulates the vertical position of the label for non expanded cells
not.expanded.label.hjust	Numeric. Regulates the horizontal position of the label for non expanded cells
total.label.vjust	Numeric. Regulates the vertical position of the center label
total.label.hjust	Numeric. Regulates the horizontal position of the center label
expanded.colors	Character vector. Colors to use for expanded clones. Should be more than 3 for better visibility. Defaults to a "darkorchid3"-based palette.
non.expanded.color	Character. Color to use for non expanded clones. Defaults to "black"

Value

Returns a list of circular plots showing proportions of expanded clones and non-expanded clones. One plot is generated for each sample in the sample_id column

Examples

```
VDJ_clonal_donut(VDJ = Platypus::small_vgm[[1]])
```

`VDJ_clonal_expansion` *Clonal frequency plot displaying clonal expansion for either T and B cells with Platypus v3 input. Only available for Platypus "v3" available. For v2 plotting of B cell clonotype expansion and isotypes please refer to `VDJ_isotypes_per_clone`.*

Description

Clonal frequency plot displaying clonal expansion for either T and B cells with Platypus v3 input. Only available for Platypus "v3" available. For v2 plotting of B cell clonotype expansion and isotypes please refer to `VDJ_isotypes_per_clone`.

Usage

```
VDJ_clonal_expansion(
  VDJ,
  celltype,
  clones,
  subtypes,
  isotypes.to.plot,
  species,
  treat.incomplete.clones,
  treat.incomplete.cells,
  group.by,
  color.by
)
```

Arguments

<code>VDJ</code>	VDJ dataframe generated using the <code>VDJ_GEX_matrix</code> function (<code>VDJ_GEX_matrix.output[[1]]</code>)
<code>celltype</code>	Character. Either "Tcells" or "Bcells". If set to Tcells bars will not be colored by default and the parameters <code>treat_incomplete_cells</code> , <code>treat_incomplete_clones</code> , <code>subtypes</code> and <code>species</code> are ignored. The <code>color.by</code> and <code>group.by</code> arguments work identically for both celltypes. If none provided it will detect this param from the <code>celltype</code> column.
<code>clones</code>	numeric value indicating the number of clones to be considered for the clonal expansion plot. Default value is 50. For a standard plot more than 50 is discouraged. When showing only one - possibly rare - isotype via <code>isotypes.to.plot</code> it may be useful to set this number higher (e.g. 100-200)
<code>subtypes</code>	Logical indicating whether to display isotype subtypes or not.

<code>isotypes.to.plot</code>	Character vector. Defaults to "all". This can be set to any number of specific Isotypes, that are to be shown exclusively. For example, to show only clones containing IgG, input "IGHG". If only wanting to check clones with IgA and IgD input c("IGHA","IGHD"). Works equally if subtypes are set to TRUE. Is ignored if color.by is not set to "isotype"
<code>species</code>	Character indicating whether the samples are from "Mouse" or "Human". Default is "Human".
<code>treat.incomplete.clones</code>	Character indicating how to proceed with clonotypes lacking a VDJC (in other words, no cell within the clonotype has a VDJC). "exclude" removes these clonotypes from the analysis. "include" keeps these clonotypes in the analysis. In the plot they will appear as having an unknown isotype.
<code>treat.incomplete.cells</code>	Character indicating how to proceed with cells assigned to a clonotype but missing a VDJC. "proportional" to fill in the VDJ isotype according to the proportions present in of clonotype (in case present proportions are not replicable in the total number of cells e.g. 1/3 in 10 cells, values are rounded to the next full integer and if the new counts exceed the total number of cells, 1 is subtracted from the isotype of highest frequency. If the number is below the number of cell, 1 is added to the isotype with lowest frequency to preserve diversity), "exclude" to exclude them from analysis and rank clonotypes only by the number of cells with a heavy chain. This ranking may deviate from the frequency column in the clonotype table. CAVE: if <code>treat_incomplete_cells</code> is set to "exclude", clonotypes lacking a VDJC entirely will be removed from the analysis. This results in a similar but not identical output as when <code>treat_incomplete_clones</code> is set to true. The two parameters are thereby non-redundant.
<code>group.by</code>	Character. Defaults to "sample_id". Column name of VDJ to split VDJ by. For each unique entry in that column a plot will be generated. Therefore plots can be generated by <code>sample_id</code> , <code>group_id</code> or any other metadata item. To get plots for the whole repertoire set to "none"
<code>color.by</code>	Character. Defaults to "isotype". If set to "isotype" bars are colored by the respective IgH chain or in grey for T cells. This can alternatively be set to any column name of the VDJ. This allows coloring clones by their <code>V_gene</code> usage or by GEX clusters

Value

Returns a nested list. `out[[1]]` are plots `out[[2]]` are raw datatables containing also barcode and CDR3 information

Examples

```
#Standard B cell plot for platypus version v3
#Will generate one plot per sample (from sample_id column)
clonal_out <- VDJ_clonal_expansion(VDJ = Platypus::small_vgm[[1]],
  celltype = "Bcells", clones = 30, subtypes = FALSE, species = "Mouse"
  , treat.incomplete.clones = "exclude"
  , treat.incomplete.cells = "proportional")
```

```

#Regrouped and recolored plot in v3
#Will generate a plot for each sample.
#Bars are filled by the sample with the highest proportion of cells in a given clonotype
clonal_out <- VDJ_clonal_expansion(VDJ = Platypus::small_vgm[[1]]
, celltype = "Bcells", clones = 30,subtypes = FALSE, species = "Mouse"
,treat.incomplete.clones = "exclude"
,treat.incomplete.cells = "proportional"
,color.by = "seurat_clusters") #change grouping with group.by = "column name"
clonal_out[[1]] #list of plots
clonal_out[[2]] #list of source dataframes

#T cell plot with recoloring by vgene
#VDJ_clonal_expansion(VDJ = VDJ.GEX.matrix.out[[1]]
#,celltype = "Tcells", clones = 30, group.by = "sample_id"
#,color_by = "VDJ_vgene")

#Plotting only IgD clones. Increased the value for clones to scan more of the dataset
#VDJ_clonal_expansion(VDJ = VDJ_comb[[1]]
#,celltype = "Bcells", clones = 150,subtypes = FALSE
#,species = "Mouse",treat.incomplete.clones = "include"
#,treat.incomplete.cells = "proportional", isotypes.to.plot = "IGHD")

#Plotting only clones containing cells with the IGHG2c isotype (For murine data only!)
#VDJ_clonal_expansion(VDJ = VDJ_comb[[1]]
#,celltype = "Bcells", clones = 150,subtypes = TRUE, species = "Mouse"
#,treat.incomplete.clones = "proportional"
#,treat.incomplete.cells = "proportional", isotypes.to.plot = "IGHG2c")

```

VDJ_clonal_lineages *Only Platypus V2 Organizes and extracts full-length sequences for clonal lineage inference. The output sequence can either contain the germline sequence as determined by cellranger or can just contain the sequences contained in each clonal family.*

Description

Only Platypus V2 Organizes and extracts full-length sequences for clonal lineage inference. The output sequence can either contain the germline sequence as determined by cellranger or can just contain the sequences contained in each clonal family.

Usage

```

VDJ_clonal_lineages(
  VDJ,
  VDJ_extract_germline.output,
  as.nucleotide,
  with.germline,
  platypus.version
)

```

Arguments

VDJ	For platypus v2 the output of the call_MIXCR function containing the full-length VDJRegion sequences. For v3 the VDJ matrix output of the VDJ_GEX_matrix function ran with trim.and.align = TRUE. (VDJ_GEX_matrix.output[[1]])
VDJ_extract_germline.output	The output from the VDJ_extract_germline function. This should have the germline information. This needs to be supplied if the with.germline argument is set to true.
as.nucleotide	Logical determining whether the full-length VDJRegion sequence should use nucleotide sequence. TRUE indicates nucleotide sequences and FALSE will extract amino acid sequences.
with.germline	Logical determining whether the germline sequence as determined by cellranger should be included in the output list of sequences. If so, the germline will be added to the last row of each dataframe object.
platypus.version	Default is "v3".

Value

returns a list containing the sequences for each clonal family as determined by the input clonotyping strategy to call_MIXCR and VDJ_extract_germline. The outer list corresponds to distinct repertoires supplied to the call_MIXCR function (e.g. VDJ.clonal.lineage.output[[i]][[j]] will contain a dataframe of the j'th clone in the i'th repertoire)

Examples

```
## Not run:
clonal_lineages <- VDJ_clonal_lineages(VDJ=call_MIXCR_output,
VDJ_extract_germline.output=VDJ_extract_germline_output,as.nucleotide=F,with.germline=T)

## End(Not run)
```

VDJ_clonotype	<i>Returns a list of clonotype dataframes following additional clonotyping. This function works best following filtering to ensure that each clone only has one heavy chain and one light chain.</i>
---------------	--

Description

Returns a list of clonotype dataframes following additional clonotyping. This function works best following filtering to ensure that each clone only has one heavy chain and one light chain.

Usage

```
VDJ_clonotype(
  VDJ,
  clone.strategy,
  homology.threshold,
  hierarchical,
  VDJ.VJ.1chain,
  global.clonotype,
  output.format,
  platypus.version
)
```

Arguments

- VDJ** For platypus v2 output from VDJ_analyze function. This should be a list of clonotype dataframes, with each list element corresponding to a single VDJ repertoire. For platypus v3 VDJ output from the VDJ_GEX_matrix function (VDJ_GEX_matrix.output[[1]])
- clone.strategy** (Updated keywords, previous format is also functional) String describing the clonotyping strategy. Possible options include 'cdr3.nt', 'cdr3.aa', 'VDJJ.VJJ', 'VDJJ.VJJ.cdr3lengths', 'VDJJ.VJJ.cdr3length.VDJcdr3homology', 'cdr3.homology', or 'VDJcdr3.homology'. 'cdr3.aa' will convert the default cell ranger clonotyping to amino acid based. 'Hvj.Lvj' groups B cells with identical germline genes (V and J segments for both heavy chain and light chain. Those arguments including 'CDR3length' will group all sequences with identical CDRH3 and CDRL3 sequence lengths. Those arguments including 'CDR3homology' will additionally impose a homology requirement for CDRH3 and CDRL3 sequences. 'CDR3homology', or 'CDRH3homology' will group sequences based on homology only (either of the whole CDR3 sequence or of the CDRH3 sequence respectively). All homology calculations are performed on the amino acid level.
- homology.threshold** Numeric value between 0 and 1 corresponding to the homology threshold for the clone.strategy arguments that require a homology threshold. Default value is set to 70 percent sequence homology. For 70 percent homology, 0.3 should be supplied as input.
- hierarchical** Boolean. Defaults to FALSE. This is an extension specifically for cells with aberrant numbers of chains (i.e. 0VDJ 1VJ, 1VDJ 0VJ, 0VDJ 2VJ, 2VDJ 0VJ). Cells with 2VDJ 2VJ are filtered out as these are most likely doublets. Aberrant cells are clonotyped hierarchically in post, following this procedure: 1. define clonotypes classically with all cells containing exactly 1VDJ 1VJ chains. 2. For cells with only a single chain (either VDJ or VJ), check if any clone exists, which matches the clonotyping criteria for this chain. If true, add this cell to that clone. If false, create a new clone containing that cell. In case that more than 1 existing clone matches the aberrant cell, the cell is assigned to the most frequent existing clone. Two reasons are behind this decision: 2.1. The aberrant cell is numerically more likely to be a part of the more frequent existing clone. 2.2 In case of a wrong assignment, the effect of the error is lower, if an already expanded

	clone is increase by one count, rather than a existing non-expanded clone being assigned a second entry and thereby resulting as expanded. 3. For cells with 3 chains, verify the clonotyping criteria on both combinations of chains (i.e. VDJ1 - VJ1, VDJ2-VJ1 in case of a cell with 2VDJ 1VJ).
VDJ.VJ.1chain	Logical specifying whether cells with multiple VDJ and VJ chains should be removed from the clonotyping. Can be either T or F for those definitions not requiring germline genes or homology thresholds, as calculating the later is difficult when multiple chains are present.
global.clonotype	Logical specifying whether clonotyping should occur across samples or only within a single sample.
output.format	String specifies function output format. Options are "vgm" (default), "dataframe.per.sample", "clone.level.dataframes", or "phylo.dataframe". "vgm" will update the existing \$clonotype_id column of the input vgm, which is the output from VDJ_GEX_matrix. "dataframe.per.sample" will return a list of VDJ dataframes, where each dataframe contains the cell-level information for a given sample. "clone.level.dataframes" will convert the per.cell matrix to a clonal dataframe, in which cells of the same clone will be merged into a single row. "dataframe.per.clone" will generate nested lists of dataframes, where each dataframe contains cell-level information of a given clone.
platypus.version	Default is "v3". To use the output of VDJ_GEX_matrix function, one should change this argument to "v3".

Value

Returns a list of clonotype dataframes where each list element matches the repertoire index in the input clonotype.list object. The dataframes will be updated with clonal frequencies based on the new clonotyping definition.

Examples

```
reclonotyped_vgm <- VDJ_clonotype(VDJ=Platypus::small_vgm[[1]],
  clone.strategy="VDJJ.VJJ",
  homology.threshold=".3", platypus.version = "v3")
```

VDJ_clonotype_clusters_circos

Makes a Circos plot from the VDJ_GEX_integrate output. Connects the clonotypes with the corresponding clusters.

Description

Makes a Circos plot from the VDJ_GEX_integrate output. Connects the clonotypes with the corresponding clusters.

Usage

```
VDJ_clonotype_clusters_circos(
  VDJ,
  topX,
  label.threshold,
  axis,
  c.count,
  n_cluster,
  platypus.version
)
```

Arguments

VDJ	The output of the VDJ_GEX_integrate function (Platypus platypus.version v2). A list of data frames for each sample containing the clonotype information and cluster membership information. For Platypus platypus.version v3, the VDJ output of the VDJ_GEX_matrix function (VDJ_GEX_matrix.output[[1]]) has to be supplied.
topX	Plots only the top X most expanded clonotypes. By default all clonotypes are shown.
label.threshold	Minimal amount of clonotypes per gene necessary to add a gene label to the sector. Default: 0.
axis	Character. Defaults to "max". Passed to VDJ_circos
c.count	Show clonotype or cell count on Circos plot. Default = T.
n_cluster	Integer. No default.
platypus.version	Which platypus.version of platypus is being used. Default = "v3".

Value

Returns list of plots. The first n elements contain the circos plot of the n datasets from the VDJ.analyze function. The n+1 element contains a list of the n adjacency matrices for each dataset.

Examples

```
#Platypus version 3
#prepare the small toy dataset
small_vgm <- Platypus::small_vgm
small_vgm[[1]]$clonotype_id_10x <- "clonotype1"
small_vgm[[1]]$clonotype_frequency <- nrow(small_vgm[[1]])
VDJ_clonotype_clusters_circos(small_vgm[[1]], topX=100, label.threshold=5
, platypus.version = "v3", n_cluster = 2)
```

VDJ_contigs_to_vgm	<i>Formats "VDJ_contigs_annotations.csv" files from cell ranger to match the VDJ_GEX_matrix output using only cells with 1VDJ and 1VJ chain</i>
--------------------	---

Description

Formats "VDJ_contigs_annotations.csv" files from cell ranger to match the VDJ_GEX_matrix output using only cells with 1VDJ and 1VJ chain

Usage

```
VDJ_contigs_to_vgm(directory, sample.names, platypus.version)
```

Arguments

directory	list containing paths to the "filtered_contig_annotations.csv" files from cell ranger.
sample.names	vector specifying sample names.
platypus.version	Function based on VGM object from V3, no need to set this parameter.

Value

data frame with column names that match the VDJ_GEX_matrix output. Can be appended to the VDJ_GEX_matrix output

Examples

```
## Not run:
directory.list <- list()
directory.list[[1]] <- c("~/Dataset_1/filtered_contig_annotations.csv")
directory.list[[2]] <- c("~/Dataset_1/filtered_contig_annotations.csv")
filtered_contig_vgm <- VDJ_contigs_to_vgm(directory = directory.list, sample.names = c(s3,s4))

## End(Not run)
```

VDJ_diversity	<i>Calculates and plots common diversity and overlap measures for repertoires and alike. Require the vegan package</i>
---------------	--

Description

Calculates and plots common diversity and overlap measures for repertoires and alike. Require the vegan package

Usage

```

VDJ_diversity(
  VDJ,
  feature.columns,
  grouping.column,
  metric,
  subsample.to.same.n,
  pvalues.label.size,
  axis.label.size,
  platypus.version
)

```

Arguments

VDJ	VDJ dataframe output from either the VDJ_analyse (platypus.version = "v2") or from the VDJ_GEX_matrix function (platypus.version = "v3")(VDJ_GEX_matrix.output[[1]])
feature.columns	Character vector. One or more column names from the VDJ of which diversity or overlap metrics are calculated. If more than one column is provided (e.g. c("VDJ_cdr3s_aa", "VJ_cdr3s_aa")) these columns will be pasted together before metric calculation. Defaults to "CDRH3_aa" if platypus.version == "v2" and "VDJ_cdr3s_aa" if platypus.version == "v3".
grouping.column	Character. Column name of a column to group metrics by. This could be "sample_id" to calculate the metric for each sample. This column is required if metric = "simpson". If so, the simpson overlap index will be calculated pairwise for all combinations of elements in the grouping.column. Defaults to "none".
metric	Character. Diversity or overlap metric to calculate. Can be c("richness", "bergerparker", "simpson", "ginisimpson", "shannon", "shannonevenness", "jaccard"). Defaults to "shannon". If jaccard is selected, a heatmap with the pairwise comparisons between all groups is returned. If any of the others is selected, a dotplot is returned
subsample.to.same.n	Boolean defaults to TRUE. Whether to subsample larger groups down to the size of the smallest group
pvalues.label.size	Numeric. Only used if overlap indices are calculated. Defaults to 4. Is passed on to ggplot theme
axis.label.size	Numeric. Only used if overlap indices are calculated. Defaults to 12. Is passed on to ggplot theme
platypus.version	Version of platypus to use. Defaults to "v3". If an output of the VDJ_analyse function is supplied, set to "v2". If an output of the VDJ_GEX_matrix function is supplied set to "v3"

Value

Returns a ggplot with the calculated metric for each group (if provided). Data is accessible via `ggplot.output$data`

Examples

```
#Calculate shannon index for VDJ CDR3s by sample
plot <- VDJ_diversity(VDJ = Platypus::small_vgm[[1]], platypus.version = "v3"
,feature.columns = c("VDJ_cdr3s_aa"), grouping.column = "sample_id"
,metric = "shannon")
#For raw values use
plot$data

#Calculate Gini-simpson and Simpson index for VDJ and VJ CDR3s by sample
VDJ_diversity(VDJ = Platypus::small_vgm[[1]], platypus.version = "v3"
,feature.columns = c("VDJ_cdr3s_aa","VJ_cdr3s_aa"), grouping.column = "sample_id"
,metric = c("ginisimpson"))

#Calculate Jaccard index of J gene usage between two samples
VDJ_diversity(VDJ = Platypus::small_vgm[[1]], platypus.version = "v3"
,feature.columns = c("VDJ_jgene"), grouping.column = "sample_id"
,metric = "jaccard")
```

VDJ_dublets	<i>Only Platypus v2 Produces a matrix indicating either the number of cells or clones which contain multiple heavy or light chains (or alpha/beta in the case of T cells).</i>
-------------	--

Description

Only Platypus v2 Produces a matrix indicating either the number of cells or clones which contain multiple heavy or light chains (or alpha/beta in the case of T cells).

Usage

```
VDJ_dublets(clonotype.list, clone.level)
```

Arguments

`clonotype.list` Output from VDJ_analyze function. This should be a list of clonotype dataframes, with each list element corresponding to a single VDJ repertoire.

`clone.level` Logical indicating whether the matrix should display information on the clone level. TRUE will result in matrices containing information about the number of chains on the clonal level. FALSE will result in matrices depicting the number of cells.

Value

Returns a list of matrices containing the number of heavy/light chains per either cell or clone depending on the clone.level parameter. This can then be supplied to heatmap functions directly. Each list element corresponds to each of the input list elements of clonotypes.

Examples

```
## Not run:
example.vdj.analyze <- VDJ_dublets(clonotype.list = "VDJ.analyze.output", clone.level=T)

## End(Not run)
```

VDJ_extract_germline *Only Platypus v2 Extracts the full-length germline sequence as determined by cellranger. This function returns an object that now contains the reference germline for each of the clones. If multiple clones (as determined by cellranger) have been merged using the VDJ_clonotype function then these sequences may have distinct germline sequences despite being in the same clonal family (nested list). This is particularly possible when homology thresholds were used to determine the clonotypes.*

Description

Only Platypus v2 Extracts the full-length germline sequence as determined by cellranger. This function returns an object that now contains the reference germline for each of the clones. If multiple clones (as determined by cellranger) have been merged using the VDJ_clonotype function then these sequences may have distinct germline sequences despite being in the same clonal family (nested list). This is particularly possible when homology thresholds were used to determine the clonotypes.

Usage

```
VDJ_extract_germline(
  VDJ.per.clone,
  mixcr.directory,
  extract.VDJRegion,
  species
)
```

Arguments

VDJ.per.clone The output from the VDJ_per_clone function. This object should have information regarding the contigs and clonotype_ids for each cell.

mixcr.directory The directory containing an executable version of MiXCR. This must be downloaded separately and is under a separate license.

extract.VDJRegion
 Default is TRUE. Future iterations will allow for distinct gene regions to be extracted.

species
 Either "mus" or "hsa" for mouse and human respectively. Default is set to mouse.

Value

Returns a dataframe containing repertoire information, such as isotype, CDR sequences, mean number of UMIs. This output can be supplied to further packages VDJ_extract_sequences and VDJ_GEX_integrate

Examples

```
## Not run:
VDJ_extract_germline(VDJ.per.clone=VDJ.per.clone.output
,mixcr.directory=~Downloads/mixcr-3.0.12/mixcr"
,extract.VDJRegion=T,species = "mmu")

## End(Not run)
```

VDJ_GEX_clonal_lineage_clusters

only Platypus v2 Integrates the transcriptional cluster information into the clonal lineages. This requires that automate_GEX, VDJ_clonal_lineages, and VDJ_GEX_integrate have already been ran. The transcriptional cluster will be added to the end of the Name for each sequence.

Description

only Platypus v2 Integrates the transcriptional cluster information into the clonal lineages. This requires that automate_GEX, VDJ_clonal_lineages, and VDJ_GEX_integrate have already been ran. The transcriptional cluster will be added to the end of the Name for each sequence.

Usage

```
VDJ_GEX_clonal_lineage_clusters(
  VDJ_GEX_integrate.output,
  VDJ_clonal_lineages.output
)
```

Arguments

VDJ_GEX_integrate.output

The output from the VDJ_GEX_integrate function that is performed on the VDJ_per_clone level. This involves a nested list where the outer list corresponds to the repertoire and inner lists correspond to specific clones based on the clonotyping strategy.

VDJ_clonal_lineages.output

Output from VDJ_clonal_lineages. This should be nested list, with the outer list element corresponding to the individual repertoire and the inner list corresponding to individual clonal lineages based on the initial clonotyping strategy in the form of a dataframe with either Seq or Name. The Name currently contains the barcode following the last "_".

Value

a nested list in the identical format to the VDJ_clonal_lineages.output but the name of each sequence will have been changed to include the transcriptional cluster corresponding to that barcode from the GEX library. This requires first running the

Examples

```
## Not run:
clonal_lineages <- VDJ_clonal_lineages(call_MIXCR.output=call_MIXCR_output
, VDJ_extract_germline.output=VDJ_extract_germline_output
,as.nucleotide=FALSE,with.germline=TRUE)

## End(Not run)
```

VDJ_GEX_expansion	<i>only Platypus v2 Integrates VDJ and gene expression libraries by providing cluster membership seq_per_vdj object. Output will plot which transcriptional cluster (GEX) that the cells of a given clonotype are found in.</i>
-------------------	---

Description

only Platypus v2 Integrates VDJ and gene expression libraries by providing cluster membership seq_per_vdj object. Output will plot which transcriptional cluster (GEX) that the cells of a given clonotype are found in.

Usage

```
VDJ_GEX_expansion(
  GEX.list,
  VDJ.GEX.integrate.list,
  highlight.isotype,
  highlight.number
)
```

Arguments

- `GEX.list` The output of the `automate_GEX` function.
- `VDJ.GEX.integrate.list` Output from `VDJ_GEX_integrate` function. This object needs to have the GEX and VDJ information combined and integrated. This should be on the CLONAL level from the `VDJ_GEX_integrate` function.
- `highlight.isotype` (Optional) isotype to plot, choose between ["None", "A", "E", "M", "G", "G1", "G2A", "G2B", "G2C", "G3"]. Default is None.
- `highlight.number` A vector corresponding to the rank of the clones that should be specified. Default is set to "20", which will present the cluster distribution for the top 20 clones.

Value

ggplot2 plot that breaks down clonotype membership per cluster for the specified input clones.

Examples

```
## Not run:
vdj.gex.expansion <- VDJ_GEX_expansion(GEX.list=GEX.list.output[[1]]
,VDJ.GEX.integrate.list=vdj.gex.integrate.output
,highlight.isotype = "None",highlight.number=1:20)

## End(Not run)
```

`VDJ_GEX_integrate` *only Platypus v2 Integrates VDJ and gene expression libraries by providing cluster membership seq_per_vdj object and the index of the cell in the Seurat RNA-seq object.*

Description

only Platypus v2 Integrates VDJ and gene expression libraries by providing cluster membership `seq_per_vdj` object and the index of the cell in the Seurat RNA-seq object.

Usage

```
VDJ_GEX_integrate(GEX.object, clonotype.list, VDJ.per.clone, clonotype.level)
```

Arguments

- `GEX.object` A single seurat object from `automate_GEX` function. This will likely be supplied as `automate_GEX.output[[1]]`.
- `clonotype.list` Output from either `VDJ_analyze` or `VDJ_clonotype` functions. This list should correspond to a single `GEX.list` object, in which each list element in `clonotype.list` is found in the `GEX.object`. Furthermore, these repertoires should be found in the `automate_GEX` library.

VDJ.per.clone Output from the VDJ_per_clone function. Each element in the list should be found in the output from the automate_GEX function.

clonotype.level

Logical specifying whether the integration should occur on the cellular level (VDJ_per_clone) or on the clonotype level (e.g. output from VDJ_analyze or VDJ_clonotype). TRUE specifies that the clonotype level will be selected - e.g. the clonotype.list object will now contain information from the GEX object regarding clonal membership.

Value

Returns a nested list containing information corresponding to either the clonal level or the sequence level, depending on the input argument "clonotype.level". This function essentially will update the output of the analyze_VDJ or the VDJ_per_clone functions.

Examples

```
## Not run:
testing_integrate <- VDJ_GEX_integrate(GEX.object = automate.gex.output[[1]]
,clonotype.list = VDJ.analyze.output
,VDJ.per.clone = VDJ.per.clone.output,clonotype.level = TRUE)

## End(Not run)
```

VDJ_GEX_matrix	<i>Processes both raw VDJ and GEX Cellranger output to compile a single cell level table containing all available information for each cell. If using Feature Barcodes please note the [FB] paragraph in the description and all "FB." parameters</i>
----------------	---

Description

This function is designed as a common input to the Platypus pipeline. Integration of datasets as well as VDJ and GEX information is done here. Please check the Platypus V3 vignette for a detailed walkthrough of the output structure. In short: output[[1]] = VDJ table, output[[2]] = GEX Seurat object and output[[3]] = statistics # [FB] Feature barcode (FB) technology is getting increasingly popular, which is why Platypus V3 fully supports their use as sample delimiters. As of V3, Platypus does not support Cite-seq data natively, also the VDJ_GEX_matrix function is technically capable of loading a Cite-seq matrix and integrating it with VDJ. For details on how to process sequencing data with FB data and how to supply this information to the VDJ_GEX_matrix function, please consult the dedicated vignette on FB data.

Usage

```
VDJ_GEX_matrix(
  VDJ.out.directory.list,
  GEX.out.directory.list,
  FB.out.directory.list,
```

```

Data.in,
Seurat.in,
VDJ.combine,
GEX.integrate,
integrate.GEX.to.VDJ,
integrate.VDJ.to.GEX,
exclude.GEX.not.in.VDJ,
filter.overlapping.barcodes.GEX,
filter.overlapping.barcodes.VDJ,
exclude.on.cell.state.markers,
get.VDJ.stats,
numcores,
trim.and.align,
select.excess.chains.by.umi.count,
gap.opening.cost,
gap.extension.cost,
parallel.processing,
integration.method,
VDJ.gene.filter,
mito.filter,
norm.scale.factor,
n.feature.rna,
n.count.rna.min,
n.count.rna.max,
n.variable.features,
cluster.resolution,
neighbor.dim,
mds.dim,
subsample.barcodes,
FB.count.threshold,
FB.ratio.threshold,
group.id,
verbose
)

```

Arguments

VDJ.out.directory.list

List containing paths to VDJ output directories from cell ranger. This pipeline assumes that the output file names have not been changed from the default 10x settings in the /outs/ folder. This is compatible with B and T cell repertoires. ! Necessary files within this folder: filtered_contig_annotations.csv, clonotypes.csv, concat_ref.fasta, all_contig_annotations.csv (only if trim.and.align == T) and metrics_summary.csv (Optional, will be appended to stats table if get.VDJ.stats == T)

GEX.out.directory.list

List containing paths the outs/ directory of each sample or directly the raw or filtered_feature_bc_matrix folder. Order of list items must be the same as for VDJ.

<code>FB.out.directory.list</code>	[FB] List of paths pointing at the outs/ directory of output from the Cellranger counts function which contain Feature barcode counts. ! Single list elements can be a path or "PLACEHOLDER", if the corresponding input in the VDJ or GEX path does not have any adjunct FB data. This is only the case when integrating two datasets of which only one has FB data. See examples for details. Any input will overwrite potential FB data loaded from the GEX input directories. This may be important, if wanting to input unfiltered FB data that will cover also cells in VDJ not present in GEX.
<code>Data.in</code>	Input for R objects from either the <code>PlatypusDB_load_from_disk</code> or the <code>PlatypusDB_fetch</code> function. If provided, input directories should not be specified. If you wish to integrate local and downloaded data, please load them via <code>load_from_disk</code> and <code>fetch</code> and provide as a list (e.g. <code>Data.in = list(load_from_disk.output, fetch.output)</code>)
<code>Seurat.in</code>	Alternative to <code>GEX.out.directory.list</code> . List of processed (!) <code>seurat</code> objects. Length of the list can either be 1 if <code>VDJ.integrate = TRUE</code> or equal to the length of <code>VDJ.out.directory.list</code> if <code>VDJ.integrate = FALSE</code> . In metadata the column <code>sample_id</code> and <code>group_id</code> must be present. <code>sample_id</code> must contain ids in the format "s1", "s2" ... "sn" and must be matching the order of <code>VDJ.out.directory.list</code> . No processing (i.e. data normalisation and integration) will be performed on these objects. They will be returned as part of the VGM and with additional VDJ data if <code>integrate.VDJ.to.GEX = T</code> . Filtering parameters such as overlapping barcodes, <code>exclude.GEX.not.in.VDJ</code> and <code>exclude.on.cell.state.markers</code> will be applied to the <code>Seurat.in</code> GEX object(s).
<code>VDJ.combine</code>	Boolean. Defaults to <code>TRUE</code> . Whether to integrate repertoires. A sample identifier will be appended to each barcode both in GEX as well as in VDJ. Recommended for all later functions
<code>GEX.integrate</code>	Boolean. Defaults to <code>TRUE</code> . Whether to integrate GEX data. Default settings use the <code>seurat scale.data</code> option to integrate datasets. Sample identifiers will be appended to each barcode both in GEX and VDJ This is helpful when analysing different samples from the same organ or tissue, while it may be problematic when analysing different tissues.
<code>integrate.GEX.to.VDJ</code>	Boolean. defaults to <code>TRUE</code> . Whether to integrate GEX metadata (not raw counts) into the VDJ output dataframe ! Only possible, if <code>GEX.integrate</code> and <code>VDJ.combine</code> are either both <code>FALSE</code> or both <code>TRUE</code>
<code>integrate.VDJ.to.GEX</code>	Boolean. defaults to <code>TRUE</code> . Whether to integrate VDJ data into GEX <code>seurat</code> object as metadata. ! Only possible, if <code>GEX.integrate</code> and <code>VDJ.combine</code> are either both <code>FALSE</code> or both <code>TRUE</code>
<code>exclude.GEX.not.in.VDJ</code>	Boolean. defaults to <code>FALSE</code> . Whether to delete all GEX cell entries, for which no VDJ information is available. Dependent on data quality and sequencing depth this may reduce the GEX cell count by a significant number
<code>filter.overlapping.barcodes.GEX</code>	Boolean. defaults to <code>TRUE</code> . Whether to remove barcodes which are shared among samples in the GEX analysis. Shared barcodes normally appear at a very low rate.

- `filter.overlapping.barcodes.VDJ`
 Boolean. defaults to TRUE. Whether to remove barcodes which are shared among samples in the GEX analysis. Shared barcodes normally appear at a very low rate.
- `exclude.on.cell.state.markers`
 Character vector. If no input is provided or input is "none", no cells are excluded. Input format should follow: Character vector containing the gene names for each state. ; is used to use multiple markers within a single gene state. Different vector elements correspond to different states. Example: c("CD4+;CD44-", "CD4+;IL7R+;CD44+"). All cells which match any of the given states (in the example case any of the 2) are excluded. This is useful in case different and non lymphocyte cells were co-sequenced. It should give the option to e.g. exclude B cells in the analysis of T cells in a dataset.
- `get.VDJ.stats` Boolean. defaults to TRUE. Whether to generate general statistics table for VDJ repertoires. This is appended as element [[3]] of the output list.
- `numcores` Number of cores used for parallel processing. Defaults to number of cores available. If you want to check how many cores are available use the library Parallel and its command detectCores() (Not setting a limit here when running this function on a cluster may cause a crash)
- `trim.and.align` Boolean. Defaults to FALSE. Whether to trim VJ/VDJ seqs, align them to the 10x reference and trim the reference. This is useful to get full sequences for antibody expression or numbers of somatic hypermutations. !Setting this to TRUE significantly increases computational time
- `select.excess.chains.by.umi.count`
 Boolean. Defaults to FALSE. There are several methods of dealing with cells containing reads for more than 1VDJ and 1VJ chain. While many analyses just exclude such cells, the VGM is designed to keep these for downstream evaluation (e.g. in VDJ_clonotype). This option presents an evidenced-based way of selectively keeping only one of the present VDJ and VJ chains each. E.g. if set to TRUE from all present VDJ chains in a cell only the one with the highest UMI count is kept. (in case of same UMI counts, the chain first in the contig files is kept). Idea source: Zhang W et al. Sci Adv. 2021 (10.1126/sciadv.abf5835)
- `gap.opening.cost`
 Argument passed to Biostrings::pairwiseAlignment during alignment to reference. Defaults to 10
- `gap.extension.cost`
 Argument passed to Biostrings::pairwiseAlignment during alignment to reference. Defaults to 4
- `parallel.processing`
 Character string. Can be "parlapply" for Windows system, "mclapply" for unix and Mac systems or "none" to use a simple for loop (slow!). Default is "none" for compatibility reasons. For the parlapply option the packages parallel, doParallel and the dependency foreach are required
- `integration.method`
 String specifying which data normalization and integration pipeline should be used. Default is "scale.data", which correspondings to the ScaleData function internal to harmony package. 'anchors' scales data individually and then finds and

align cells in similar states as described here: https://satijalab.org/seurat/articles/integration_introduction.html. 'sct' specifies SCTransform from the Seurat package. "harmony" should be specified to perform harmony integration. This method requires the harmony package from bioconductor.

VDJ.gene.filter	Logical indicating if variable genes from the b cell receptor and t cell receptor should be removed from the analysis. True is highly recommended to avoid clonal families clustering together.
mito.filter	Numeric specifying which percent of genes are allowed to be composed of mitochondrial genes. This value may require visual inspection and can be specific to each sequencing experiment. Users can visualize the percentage of genes corresponding to mitochondrial genes using the function "investigate_mitochondrial_genes".
norm.scale.factor	Scaling factor for the standard Seurat pipeline. Default is set to 10000 as reported in Seurat documentation.
n.feature.rna	Numeric that specifies which cells should be filtered out due to low number of detected genes. Default is set to 0. Seurat standard pipeline uses 2000.
n.count.rna.min	Numeric that specifies which cells should be filtered out due to low RNA count. Default is set to 0. Seurat standard pipeline without VDJ information uses 200.
n.count.rna.max	Numeric that specifies which cells should be filtered out due to high RNA count. Default is set to infinity. Seurat standard pipeline without VDJ information uses 2500.
n.variable.features	Numeric specifying the number of variable features. Default set to 2000 as specified in Seurat standard pipeline.
cluster.resolution	Numeric specifying the resolution that will be supplied to Seurat's FindClusters function. Default is set to 0.5. Increasing this number will increase the number of distinct Seurat clusters. Suggested to examine multiple parameters to ensure gene signatures differentiating clusters remains constant.
neighbor.dim	Numeric vector specifying which dimensions should be supplied in the FindNeighbors function from Seurat. Default input is '1:10'.
mds.dim	Numeric vector specifying which dimensions should be supplied into dimensional reduction techniques in Seurat and Harmony. Default input is '1:10'.
subsample.barcodes	For development purposes only. If set to TRUE the function will run on 100 cells only to increase speeds of debugging
FB.count.threshold	Numeric. Defaults to 10. For description of Feature Barcode assignment see parameter FB.ratio.threshold above
FB.ratio.threshold	Numeric. Defaults to 2 Threshold for assignment of feature barcodes by counts. A feature barcode is assigned to a cell if its counts are >FB.count.threshold and if its counts are FB.ratio.threshold-times higher than the counts of the feature barcode with second most counts.

group.id vector with integers specifying the group membership. c(1,1,2,2) would specify the first two elements of the input VDJ/GEX lists are in group 1 and the third/fourth input elements will be in group 2.

verbose if TRUE prints runtime info to console. Defaults to TRUE

Value

Single cell matrix including VDJ and GEX info. Format is a list with out[[1]] = a VDJ dataframe (or list of dataframes if VDJ.combine == F, not recommended) containing also selected GEX information of integrate.GEX.to.VDJ = T. out[[2]] = GEX Seurat object with the metadata also containing GEX information if integrate.VDJ.to.GEX = T. out[[3]] = Dataframe with statistics on GEX and VDJ. out[[4]] = runtime parameters. out[[5]] = session info

Examples

```
## Not run:

#FOR EXAMPLES see Platypus vignette at https://alexgermanos.github.io/Platypus/index.html

#Run from local directory input. For run from PlatypusDB input see
#PlatypusDB vignette
VDJ.out.directory.list <- list()
VDJ.out.directory.list[[1]] <- c("~/VDJ/S1/")
VDJ.out.directory.list[[2]] <- c("~/VDJ/S2/")
GEX.out.directory.list <- list()
GEX.out.directory.list[[1]] <- c("~/GEX/S1/")
GEX.out.directory.list[[2]] <- c("~/GEX/S2/")
VGM <- VDJ_GEX_matrix(
  VDJ.out.directory.list = VDJ.out.directory.list
  ,GEX.out.directory.list = GEX.out.directory.list
  ,GEX.integrate = T
  ,VDJ.combine = T
  ,integrate.GEX.to.VDJ = T
  ,integrate.VDJ.to.GEX = T
  ,exclude.GEX.not.in.VDJ = F
  ,filter.overlapping.barcodes.GEX = F
  ,filter.overlapping.barcodes.VDJ = F
  ,get.VDJ.stats = T
  ,parallel.processing = "none"
  ,subsample.barcodes = F
  ,trim.and.align = F
  ,group.id = c(1,2))

# With Feature Barcodes
## Option 1: Cellranger multi or Cellranger count with --libraries output
VDJ.out.directory.list <- list()
VDJ.out.directory.list[[1]] <- "~/VDJ/S1/" #point to outs or per_sample_outs directory content
VDJ.out.directory.list[[2]] <- "~/VDJ/S2/"
GEX.out.directory.list <- list()
GEX.out.directory.list[[1]] <- "~/GEX/S1/"
GEX.out.directory.list[[2]] <- "~/GEX/S2/" #These directories contain two matrices (GEX and FB)
VGM <- VDJ_GEX_matrix(
```

```

VDJ.out.directory.list = VDJ.out.directory.list
,GEX.out.directory.list = GEX.out.directory.list,
FB.ratio.threshold = 2)

##Option 2: Separate input of FB data from separate Cellranger count run
VDJ.out.directory.list <- list()
VDJ.out.directory.list[[1]] <- "~/VDJ/S1/"
VDJ.out.directory.list[[2]] <- "~/VDJ/S2/"
GEX.out.directory.list <- list()
GEX.out.directory.list[[1]] <- "~/GEX/S1/"
GEX.out.directory.list[[2]] <- "~/GEX/S2/"
GEX.out.directory.list <- list()
FB.out.directory.list[[1]] <- "~/FB/S1/"
FB.out.directory.list[[2]] <- "~/FB/S1/"
VGM <- VDJ_GEX_matrix(
VDJ.out.directory.list = VDJ.out.directory.list,
GEX.out.directory.list = GEX.out.directory.list,
FB.out.directory.list = FB.out.directory.list,
FB.ratio.threshold = 2)

##Option 3: FB input for two datasets of which only one contains FB data
VDJ.out.directory.list <- list()
VDJ.out.directory.list[[1]] <- "~/study1/VDJ/S1/"
VDJ.out.directory.list[[2]] <- "~/study2/VDJ/S1/"
VDJ.out.directory.list[[3]] <- "~/study2/VDJ/S2/"
GEX.out.directory.list <- list()
GEX.out.directory.list[[1]] <- "~/study1/GEX/S1/"
GEX.out.directory.list[[2]] <- "~/study2/GEX/S1/"
GEX.out.directory.list[[2]] <- "~/study2/GEX/S2/"
GEX.out.directory.list <- list()
FB.out.directory.list[[1]] <- "PLACEHOLDER" #Study 1 does not contain FB data
FB.out.directory.list[[2]] <- "~/study2/FB/S1/"
FB.out.directory.list[[3]] <- "~/study2/FB/S2/"
VGM <- VDJ_GEX_matrix(
VDJ.out.directory.list = VDJ.out.directory.list,
GEX.out.directory.list = GEX.out.directory.list,
FB.out.directory.list = FB.out.directory.list,
FB.ratio.threshold = 2)

## End(Not run)

```

VDJ_GEX_overlay_clones

Highlights the cells belonging to any number of top clonotypes or of specifically selected clonotypes from one or more samples or groups in a GEX dimensional reduction.

Description

Highlights the cells belonging to any number of top clonotypes or of specifically selected clonotypes from one or more samples or groups in a GEX dimensional reduction.

Usage

```
VDJ_GEX_overlay_clones(
  GEX,
  reduction,
  n.clones,
  clones.to.plot,
  by.sample,
  by.other.group,
  ncol.facet,
  pt.size,
  clone.colors,
  split.plot.and.legend,
  platypus.version
)
```

Arguments

GEX	A single <code>seurat</code> object from <code>VDJ_GEX_matrix</code> , which also includes VDJ information in the metadata (set <code>integrate.VDJ.to.GEX</code> to <code>TRUE</code> in the <code>VDJ_GEX_matrix</code> function) (<code>VDJ_GEX_matrix.output[[2]]</code>)
reduction	Character. Defaults to "umap". Name of the reduction to overlay clones on. Can be "pca", "umap", "tsne"
n.clones	Integer. Defaults to 5. TO PLOT TOP N CLONES. Number of Top clones to plot. If either <code>by.sample</code> or <code>by.group</code> is <code>TRUE</code> , <code>n.clones</code> clones from each sample or group will be overlaid
clones.to.plot	Character. Alternative to <code>n.clones</code> . TO PLOT SPECIFIC CLONES. Must reference a column in the <code>GEX@meta.data</code> filled with <code>TRUE</code> and <code>FALSE</code> . Entries with <code>TRUE</code> label are plotted. Such a column may be generated using <code>GEX@metadata\$clones_to_plot_column <- GEX@metadata\$Some_cell_identifier == "Interesting"</code>
by.sample	Boolean. Defaults to <code>FALSE</code> . Whether to overlay clones by sample. If set to <code>TRUE</code> this will generate a <code>facet_wrap</code> plot with as many facets as samples.
by.other.group	Character string. Defaults to "none". Must be a valid column name of the metadata of the input <code>seurat</code> object. If so, this will generate a <code>facet_wrap</code> plot with as many facets unique entries in the specified column. This may be useful to plot cell type specific clones
ncol.facet	Integer. Defaults to 2. Number of columns in the <code>facet_wrap</code> plot if <code>by.sample</code> or <code>by.group</code> is <code>TRUE</code>
pt.size	Numeric. Defaults to 1. Size of points in <code>DimPlot</code> . Passed to <code>Seurat::DimPlot</code>
clone.colors	Character vector. Defaults to <code>rainbow(n.clones)</code> . Colors to use for individual clones. One can provide either a vector of length <code>n.clones</code> or a of length <code>Nr. of</code>

`n.clones` `*` `n.clones`. In case that a vector of length `n.clones` is provided and `by.group` or `by.sample` is `TRUE`, colors are repeated for each sample/group

`split.plot.and.legend`
 Boolean. Defaults to `FALSE`. Whether to return the plot and the legend separately as a list. This can be useful if legends get large and distort the actual plots. The packages `gridExtra` and `cowplot` are required for this. If set to `TRUE` a list is returned where `out[[1]]` is the plot which can be printed just by executing `out[[1]]`; `out[[2]]` is the legend, which can be printed either using `plot(out[[2]])` or `grid.arrange(out[[2]])`

`platypus.version`
 Character. At the moment this function runs only on the output of the `VDJ_GEX_matrix` function meaning that it is exclusively part of Platypus "v3". With further updates the functionality will be extended.

Value

A `ggplot` object or a list of a `ggplot` and a `gtable` legend (if `split.plot.and.legend` `\=` `TRUE`). Theme, colors etc. may be changed directly by adding new elements to this output (e.g. `out` `\+` `theme_minimal()`)

Examples

```
#To return a single plot with top clones across samples
overlay_clones_plot <- VDJ_GEX_overlay_clones(
  GEX = Platypus::small_vgm[[2]], reduction = "umap"
  ,n.clones = 5, by.sample = FALSE
  ,by.other.group = "none", pt.size = 1,split.plot.and.legend = FALSE)

#To return a facet plot with top clones for each sample
overlay_clones_plot <- VDJ_GEX_overlay_clones(
  GEX = Platypus::small_vgm[[2]], reduction = "umap"
  ,n.clones = 5, by.sample = TRUE, by.other.group = "none"
  ,pt.size = 1,ncol.facet = 2, split.plot.and.legend = FALSE)

#To return a facet plot and the legend separately with top clones for each group
overlay_clones_plot <- VDJ_GEX_overlay_clones(
  GEX = Platypus::small_vgm[[2]], reduction = "umap"
  ,n.clones = 5, by.sample = TRUE, by.other.group = "group_id", pt.size = 1
  ,ncol.facet = 2, split.plot.and.legend = TRUE)

#To print both:
#overlay_clones_plot[[1]] #Plot
#gridExtra::grid.arrange(overlay_clones_plot[[2]]) #Legend
#To save, ggsave() is applicable to both

#To return a single plot with selected clones
#add a clonotype_to_plot column
#GEX@meta.data$clonotype_to_plot <- GEX$VJ_vgene == "TRAV5-1"
#Column with TRUE for all clones with a particular V gene
#overlay_clones_plot <- VDJ_GEX_overlay_clones(GEX = GEX, reduction = "umap")
```

```
#, clones.to.plot = "clonotype_to_plot", by.sample = TRUE, by.other.group = "none"
#, split.plot.and.legend = FALSE, pt.size = 1.5)
```

VDJ_GEX_stats	<i>Gives stats on number and quality of reads.</i>
---------------	--

Description

Gives stats on number and quality of reads.

Usage

```
VDJ_GEX_stats(
  VDJ.out.directory,
  GEX.out.directory,
  sample.names,
  metrics10x,
  save.csv,
  filename
)
```

Arguments

VDJ.out.directory	List of paths with each element containing the path to the output of cellranger VDJ runs. This pipeline assumes that the output file names have not been changed from the default 10x settings in the /outs/ folder. This is compatible with B and T cell repertoires (both separately and simultaneously).
GEX.out.directory	OPTIONAL list of paths with each element containing the path to the output of cellranger GEX runs. This pipeline assumes that the output file names have not been changed from the default 10x settings in the /outs/ folder. This is compatible with B and T cell repertoires (both separately and simultaneously).
sample.names	OPTIONAL: an array of the same length as the input VDJ.out.directory list with custom names for each sample. If not provided samples will be numbered by processing order
metrics10x	Whether to append metrics_summary.csv information provided by Cellranger for both VDJ and GEX. Defaults to T
save.csv	Boolean. Defaults to TRUE. Whether to directly save the results as a comma delimited .csv file in the current working directory.
filename	Character ending in .csv. Filename to save .csv as.

Value

returns a single matrix where the rows are individual cells and the columns are repertoire features.

Examples

```
## Not run:
stats <- VDJ_GEX_stats(VDJ.out.directory = VDJ.out.directory.list
,GEX.out.directory = GEX.out.directory.list,sample.names = c(1:4)
,metrics10x = TRUE,save.csv = TRUE ,filename = "stats.csv")

## End(Not run)
```

VDJ_isotypes_per_clone

Only for Platypus v2 Clonal frequency plot displaying the isotype usage of each clone. ! For platypus v3 use VDJ_clonal_expansion

Description

Only for Platypus v2 Clonal frequency plot displaying the isotype usage of each clone. ! For platypus v3 use VDJ_clonal_expansion

Usage

```
VDJ_isotypes_per_clone(
  VDJ_clonotype_output,
  VDJ_per_clone_output,
  clones,
  subtypes,
  species,
  sample.names,
  treat.incomplete.clones,
  treat.incomplete.cells,
  platypus.version,
  VDJ.matrix
)
```

Arguments

VDJ_clonotype_output	list of dataframes based on the VDJ_clonotype function output.
VDJ_per_clone_output	list of dataframes based on the VDJ_per_clone function output.
clones	numeric value indicating the number of clones to be displayed on the clonal expansion plot. Can take values between 1-50. Default value is 50.
subtypes	Logical indicating whether to display isotype subtypes or not.
species	Character indicating whether the samples are from mouse or human. Default is set to human. #' @param sample.names Character vector with the same length of the VDJ.GEX.matrix.out list. If a VDJ table is provided, length of samples names must be one. These names are used as references to the output and as title for the plots

`sample.names` Vector. Names for samples in the order of the `VDJ_GEX_matrix` or the `VDJ.analyze.output`. Defaults to 1-n

`treat.incomplete.clones`
Character indicating how to proceed with clonotypes lacking a VDJC (in other words, no cell within the clonotype has a VDJC). "exclude" removes these clonotypes from the analysis. This may result in a different frequency ranking of clonotypes than in the output of the `VDJ_analyse` function with `filter.IHC.ILC = FALSE`. "include" keeps these clonotypes in the analysis. In the plot they will appear as having an unknown isotype.

`treat.incomplete.cells`
Character indicating how to proceed with cells assigned to a clonotype but missing a VDJC. "proportional" to fill in the VDJ isotype according to the proportions present in of clonotype (in case present proportions are not replicable in the total number of cells e.g. 1/3 in 10 cells, values are rounded to the next full integer and if the new counts exceed the total number of cells, 1 is subtracted from the isotype of highest frequency. If the number is below the number of cell, 1 is added to the isotype with lowest frequency to preserve diversity), "exclude" to exclude them from analysis and rank clonotypes only by the number of actual contigs of there heavy chain. This ranking may deviate from the frequency column in the clonotype table. CAVE: if `treat_incomplete_cells` is set to "exclude", clonotypes lacking a VDJC entirely will be removed from the analysis. This results in a similar but not identical output as when `treat_incomplete_clones` is set to true. The two parameters are thereby non-redundant.

`platypus.version`
Defaults to "v3". For a more flexible analysis in v3 use `VDJ_clonal_expansion()`

`VDJ.matrix` The VDJ table output of the `VDJ_GEX_matrix` function. (`VDJ_GEX_matrix.output[[1]]`)

Value

returns a list containing plots with the percentages of isotypes for each clone on the cell level.

Examples

```
## Not run:
VDJ.isotype.per.clone <- VDJ_isotypes_per_clone(
VDJ_clonotype_output = VDJ.analyze.output
,VDJ_per_clone_output = VDJ.per.clone.output, clones = 30)

## End(Not run)
```

vdj_length_prob	<p><i>vdj_length_prob</i> A list dataframe specifying lengths and probabilities of bases deleted or inserted at each junction site of VDJ recombination event.</p> <p>v3_deletion length and probability of deleted bases at 3' end of V segment</p> <p>d5_deletion length and probability of deleted bases at 5' end of D segment</p> <p>d3_deletion length and probability of deleted bases at 3' end of D segment</p> <p>j5_deletion length and probability of deleted bases at 5' end of J segment</p> <p>dj_insertion length and probability of inserted bases between D-J segment</p> <p>vj_insertion length and probability of inserted bases between V-J segment for light or alpha chains</p>
-----------------	--

Description

vdj_length_prob A list dataframe specifying lengths and probabilities of bases deleted or inserted at each junction site of VDJ recombination event.

v3_deletion length and probability of deleted bases at 3' end of V segment

d5_deletion length and probability of deleted bases at 5' end of D segment

d3_deletion length and probability of deleted bases at 3' end of D segment

j5_deletion length and probability of deleted bases at 5' end of J segment

dj_insertion length and probability of inserted bases between D-J segment

vj_insertion length and probability of inserted bases between V-J segment for light or alpha chains

Usage

```
data("vdj_length_prob")
```

Format

An object of class `list` of length 7.

VDJ_logoplot_vector *Plots a logoplot of the CDR3 aminoacid region*

Description

Plots a logoplot of the CDR3 aminoacid region

Usage

```
VDJ_logoplot_vector(cdr3.vector, length_cdr3, seq_type)
```

Arguments

cdr3.vector	A character vector of aa sequences. This is to increase flexibility of this function. Such a sequence vector may be retrieved from the VDJ_analyse function output on a clonotype level or from the VDJ_GEX_matrix function output on a per cell level. Additionally, any length of sequence may be used (e.g. HCDR3 only or H and LCDR3 pasted together)
length_cdr3	Integer or character. Defaults to "auto". Sets the length of the CDR3 regions that are selected to be plotted. If set to auto, the most frequently appearing length in the vector will be used
seq_type	passed to ggseqlogo. Can be set to "aa", "dna", "rna" or "other"

Value

Returns the logo plot.

Examples

```
VDJ_logoplot_vector(
  cdr3.vector = Platypus::small_vgm[[1]]$VDJ_cdr3s_aa
  ,length_cdr3 = "auto",seq_type = "auto")
```

VDJ_network *Creates a similarity network where clones with similar CDR3s are connected.*

Description

Creates a similarity network where clones with similar CDR3s are connected.

Usage

```
VDJ_network(
  VDJ,
  distance.cutoff,
  per.sample,
  platypus.version,
  known.binders,
  hcdr3.only
)
```

Arguments

VDJ	Either (for platypus version "v2") output from VDJ_analyze function. This should be a list of clonotype dataframes, with each list element corresponding to a single VDJ repertoire, OR (for platypus version "v3") the the VDJ matrix output of the VDJ_GEX_matrix() function (VDJ.GEX.matrix.output[[1]])
distance.cutoff	The threshold Levenshtein distance for which two nodes will be connected on the similarity network.
per.sample	logical value indicating if a single networks should be produced for each mouse.
platypus.version	Character. Defaults to "v3". Can be "v2" or "v3" dependent on the input format
known.binders	Either a character vector with cdr3s of known binders or a data frame with cdr3s in the first and the corresponding specificity in the second column. If this parameter is defined, the output will be a network with only edges between known binders and the repertoire nodes and edges between the known binders that have at least one edge to a repertoire node
hcdr3.only	logical value indicating if the network is based on heavy chain cdr3s (hcdr3.only = T) or pasted heavy and light chain cdr3s (hcdr3.only = F), works for platypus.version 3 only

Value

returns a list containing networks and network information. If per.sample is set to TRUE then the result will be a network for each repertoire. If per.sample ==F, output[[1]] <- will contain the network, output[[2]] will contain the dataframe with information on each node, such as frequency, mouse origin etc. output[[3]] will contain the connected index - these numbers indicate that the nodes are connected to at least one other node. output[[4]] contains the paired graph - so the graph where only the connected nodes are drawn.

Examples

```
#Platypus v2
#network_out <- VDJ_network(VDJ = VDJ_analyze.out[[1]],per.sample = TRUE,distance.cutoff = 2)
#Platypus v3
network_out <- VDJ_network(VDJ = Platypus::small_vgm[[1]],per.sample = FALSE,distance.cutoff = 2)
```

VDJ_overlap_heatmap *Yields overlap heatmap and datatable of features or combined features for different samples or groups*

Description

Yields overlap heatmap and datatable of features or combined features for different samples or groups

Usage

```
VDJ_overlap_heatmap(
  VDJ,
  feature.columns,
  grouping.column,
  plot.type,
  pvalues.label.size,
  axis.label.size,
  add.barcode.table
)
```

Arguments

VDJ	VDJ output of the VDJ_GEX_matrix function (VDJ_GEX_matrix.output[[1]])
feature.columns	A character array of column names of which the overlap should be displayed. The content of these columns is pasted together (separated by "/"). E.g. if the overlap in cells germline gene usage is desired, the input could be c("VDJ_jgene", "VDJ_dgene", "VDJ_vg"). These columns would be pasted and compared across the grouping variable.
grouping.column	A column which acts as a grouping variable. If repertoires are to be compared use the sample_id column.
plot.type	Character. Either "ggplot" or "pheatmap". Defaults to Pheatmap
pvalues.label.size	Numeric. Defaults to 4. Is passed on to ggplot theme
axis.label.size	Numeric. Defaults to 4. Is passed on to ggplot theme
add.barcode.table	Boolean. Defaults to T. Whether to generate a dataframe with frequencies and barcodes of cells with overlapping features. This is useful to e.g. analyze differentially expressed genes between cells of two samples or groups expressing the same VDJ or VJ chain

Value

A list of a ggplot (out[[1]]), the source table or matrix for the plot out[[2]] and a table containing additional information in case that add.barcode.table was set to TRUE (out[[3]])

Examples

```
#To test the overlap of CDR3s between multiple samples
overlap <- VDJ_overlap_heatmap(VDJ = Platypus::small_vgm[[1]]
,feature.columns = c("VDJ_cdr3s_aa"),
grouping.column = "sample_id", axis.label.size = 15
, plot.type = "ggplot")
```

VDJ_per_clone

VDJ_per_clone

Description

only Platypus v2 Analyzes and processes the repertoire sequencing data from cellranger vdj. This provides information on the single-cell level for each clone, as opposed to the output from VDJ_analyze.

Usage

```
VDJ_per_clone(
  clonotype.list,
  VDJ.out.directory,
  contig.list,
  fasta.list,
  reference.list,
  filtered.contigs,
  annotations.json,
  JSON
)
```

Arguments

`clonotype.list` Output from either VDJ_analyze or VDJ_clonotype functions. This list should correspond to a single GEX.list object, in which each list element in clonotype.list is found in the GEX.object. Furthermore, the *i*'th entry in the directory supplied to GEX.list should correspond to the *i*'th element in the clonotype.list object.

`VDJ.out.directory`

Character vector with each element containing the path to the output of cellranger vdj runs. This corresponds to the same object used for the VDJ_analyze function. Multiple repertoires to be integrated in a single transcriptome should be supplied as multiple elements of the character vector. This can be left blank if supplying the clonotypes and contig files directly as input. This pipeline assumes that the output file names have not been changed from the default 10x settings in the /outs/ folder. This is compatible with B and T cell repertoires (both separately and simultaneously).

<code>contig.list</code>	List of dataframe based on the <code>all_contigs.csv</code> file from cellranger vdj output. If 10x sequencing was not used then this object should be formatted with the same columns as the 10x object.
<code>fasta.list</code>	Contains the full-length sequence information in the same format as <code>filtered_contig.fasta</code> file from the output of cellranger.
<code>reference.list</code>	Contains the reference sequence information in the same format as <code>concat_ref.fasta</code> file from the output of cellranger.
<code>filtered.contigs</code>	Logical indicating if the filtered contigs file should be used. TRUE will read VDJ information from only the filtered output of cellranger. FALSE will read the all contigs file from cellranger. Default set to TRUE (filtered output)
<code>annotations.json</code>	Optional input from loaded <code>all_contig_annotations.json</code> . Will be read in automatically if not provided
JSON	Boolean. Defaults to FALSE. Whether to load <code>all_contig_annotations.json</code>

Value

Returns a list of dataframes containing

Examples

```
## Not run:
VDJ_per_clone_out <- VDJ_per_clone(clonotype.list = output.from.VDJ_analyze
,VDJ.out.directory = "path/to/cellranger/outs/")

## End(Not run)
```

VDJ_plot_SHM	<i>Plots for SHM based on MIXCR output generated using the VDJ_call_MIXCR function and appended to the VDJ.GEX.matrix.output</i>
--------------	--

Description

Plots for SHM based on MIXCR output generated using the `VDJ_call_MIXCR` function and appended to the `VDJ.GEX.matrix.output`

Usage

```
VDJ_plot_SHM(
  VDJ.mixcr.matrix,
  group.by,
  quantile.label,
  point.size,
  mean.line.color,
  stats.to.console,
  platypus.version
)
```

Arguments

<code>VDJ.mixcr.matrix</code>	Output dataframe from the <code>VDJ_call_MIXCR</code> function or a dataframe generated using the <code>VDJ_GEX_matrix</code> function and supplemented with MIXCR information
<code>group.by</code>	Character. Defaults to "sample_id". Column name of <code>VDJ.matrix</code> to split <code>VDJ.matrix</code> by. For each unique entry in that column a set of plots will be generated. This can be useful to plot SHM by expansion or by transcriptomics-derived clusters
<code>quantile.label</code>	Numeric. Defaults to 0.9. Which points to label in the SHM scatterplot. If set to 0.9, the top 10% of cells by SHM number will be labelled. If <code>ggrepel</code> throws a warning, concerning overlap it is recommended to attempt to label less points to avoid cluttering
<code>point.size</code>	Size of points in plots. Passed to <code>geom_jitter()</code>
<code>mean.line.color</code>	Color of mean bar in dotplots. Passed to <code>geom_errorbar()</code>
<code>stats.to.console</code>	Boolean. Defaults to FALSE. Prints basic statistics (AOV + post hoc test) to console
<code>platypus.version</code>	Character. Only "v3" available.

Value

Returns a list of ggplot objects. `out[[1]]` is a boxplot comparing SHM by `group.by`. `out[[2]]` to `out[[n]]` are plots for each group that visualize VDJ and VJ SHM distribution for each group. Data for any plot can be accessed via `out[[any]]$data`

Examples

```
#Simulating SHM data
small_vgm <- Platypus::small_vgm
small_vgm[[1]]$VDJ_SHM <- as.integer(rnorm(nrow(small_vgm[[1]]), mean = 5, sd = 3))
small_vgm[[1]]$VJ_SHM <- as.integer(rnorm(nrow(small_vgm[[1]]), mean = 5, sd = 3))

#Standard plots
SHM_plots <- VDJ_plot_SHM(VDJ = small_vgm[[1]]
, group.by = "sample_id", quantile.label = 0.9)

#Group by transcriptional cluster and label only top 1%
SHM_plots <- VDJ_plot_SHM(VDJ = small_vgm[[1]]
, group.by = "seurat_clusters", quantile.label = 0.99)
```

VDJ_reclonotype_list_arrange

Only Platypus v2 Organizes the top N genes that define each Seurat cluster and converts them into a single dataframe. This can be useful for obtaining insight into cluster-specific phenotypes.

Description

Only Platypus v2 Organizes the top N genes that define each Seurat cluster and converts them into a single dataframe. This can be useful for obtaining insight into cluster-specific phenotypes.

Usage

```
VDJ_reclonotype_list_arrange(  
  VDJ_clonotype.output,  
  VDJ_analyze.output,  
  Platypus_list.object  
)
```

Arguments

VDJ_clonotype.output

The output object from the VDJ_clonotype function. The column of the merged nucleotide clonotype IDs will be used to rearrange the new object.

VDJ_analyze.output

The output from the initial VDJ_analyze, containing clonotype information based on nucleotide sequence.

Platypus_list.object

The new list object from one of Platypus functions (for example, clonal lineages, VDJ_per_clne, etc) that should be merged based on the VDJ_clonotype output structure. nested list structure, where outer list corresponds to repertoire and the inner list corresponds to clones (on the nucleotide level).

Value

Returns a dataframe in which the top N genes defining each cluster based on differential expression are selected.

Examples

```
## Not run:  
checking_vdj_reclono <- VDJ_reclonotype_list_arrange(  
  VDJ_clonotype.output = repertoire_reclonotype  
  ,VDJ_analyze.output = repertoire_list  
  ,Platypus_list.object = repertoire_vdj_per_clone)  
  
## End(Not run)
```

VDJ_tree	<i>only Platypus v2 Produces neighbor joining phylogenetic trees from the output of VDJ_clonal_lineages</i>
----------	---

Description

only Platypus v2 Produces neighbor joining phylogenetic trees from the output of VDJ_clonal_lineages

Usage

```
VDJ_tree(
  clonal.lineages,
  with.germline,
  min.sequences,
  max.sequences,
  normalize.germline.length,
  unique.sequences
)
```

Arguments

clonal.lineages	Output from VDJ_clonal_lineages. This should be nested list, with the outer list element corresponding to the individual repertoire and the inner list corresponding to individual clonal lineages based on the initial clonotyping strategy in the form of a dataframe with either Seq or Name.
with.germline	Logical specifying if the germline should be set as outgroup. Default is set to TRUE.
min.sequences	integer value specifying the minimum number of sequences to be allowed for clonal lineages. Default is 3.
max.sequences	integer value specifying the maximum number of sequences to be allowed for clonal lineages. Default is 500
normalize.germline.length	Logical determining whether or not the branch length separating the germline from the first internal node should be normalized. Potentially useful for visualization if the remainder tips are far from the root. Default is TRUE.
unique.sequences	Logical indicating if those cells containing identical VDJRegion sequences should be merged into single nodes and have their variant added as the tip label. Default is TRUE.

Value

Returns a nested list of phylogenetic trees. The output[[i]][[j]] corresponds to the j'th clone in the i'th input repertoire. plot(output[[i]][[j]]) should display the phylogenetic tree if the ape package is loaded.

Examples

```
## Not run:
vdj.tree <- VDJ_tree(clonal.lineages = VDJ.clonal.lineage.output
,with.germline=TRUE,min.sequences = 5
,max.sequences = 30,unique.sequences = TRUE)

## End(Not run)
```

VDJ_variants_per_clone

Returns statistics and plots to examine diversity of any sequence or metadata item within clones on a by sample level or global level

Description

Returns statistics and plots to examine diversity of any sequence or metadata item within clones on a by sample level or global level

Usage

```
VDJ_variants_per_clone(
  VDJ,
  variants.of,
  clonotypes.col,
  stringDist.method,
  split.by,
  platypus.version
)
```

Arguments

VDJ	VDJ output of the VDJ_GEX_matrix (VDJ_GEX_matrix.output[[1]]). VDJ matrix supplemented with with MIXCR information is also valid
variants.of	Character vector. Defaults to c("VDJ_cdr3s_aa", "VJ_cdr3s_aa"). Column name(s) of VDJ to examine variants of. If more than one name is given, these columns will be pasted together. The default will therefore return statistics on the number of variants of VDJ and VJ cdr3s in every clone
clonotypes.col	Column name of the VDJ column containing clonotype information. Defaults to "clonotype_id_10x". This is useful if alternative clonotyping strategies have been used and are stored in other columns
stringDist.method	Character. Passed to Biostrings::strinDist. Method to calculate distance between variants of a clone. Defaults to "levenshtein". Other options are "hamming", "quality". If "hamming" variants of a clone will be shortened from the end to the shortest variant to make all input sequences the same length.

`split.by` Character. Defaults to "sample_id". Column name of VDJ to split the analysis by. This is necessary, if clonotyping was done on a per sample level (e.g. "clonotype1" in sample 1 is not the same sequence as "clonotype1" in sample 2). If clonotyping was done across samples and no splitting is necessary input "none"

`platypus.version` Character. Only "v3" available.

Value

Returns a list of dataframes. Each dataframe contains the statistics of one `split.by` element (by default: one sample)

Examples

```
variants_per_clone <- VDJ_variants_per_clone(VDJ = Platypus::small_vgm[[1]]
,variants.of = c("VDJ_cdr3s_aa", "VJ_cdr3s_aa"),
stringDist.method = "levenshtein", split.by = "sample_id")
```

VDJ_Vgene_usage	<i>Produces a matrix counting the number of occurrences for each VDJ and VJ Vgene combinations for each list entry in VDJ.clonotype.output or for each sample_id in VDJ.matrix</i>
-----------------	--

Description

Produces a matrix counting the number of occurrences for each VDJ and VJ Vgene combinations for each list entry in VDJ.clonotype.output or for each sample_id in VDJ.matrix

Usage

```
VDJ_Vgene_usage(VDJ, platypus.version)
```

Arguments

`VDJ` For `platypus.version = "v2"` output from `VDJ_analyze` function. This should be a list of clonotype dataframes, with each list element corresponding to a single VDJ repertoire. For `platypus.version = "v3"` output VDJ dataframe from `VDJ_GEX_matrix` function (`VDJ_GEX_matrix.output[[1]]`)

`platypus.version` Character. Defaults to "v3". Can be "v2" or "v3" dependent on the input format

Value

Returns a list of matrices containing the number of Vgene heavy/light chain combinations per repertoire.

Examples

```
example.vdj.vgene_usage <- VDJ_Vgene_usage(VDJ =
  Platypus::small_vgm[[1]], platypus.version = "v3")
```

VDJ_Vgene_usage_barplot

Produces a barplot with the most frequently used IgH and IgK/L Vgenes.

Description

Produces a barplot with the most frequently used IgH and IgK/L Vgenes.

Usage

```
VDJ_Vgene_usage_barplot(
  VDJ,
  HC.gene.number,
  LC.Vgene,
  LC.gene.number,
  platypus.version
)
```

Arguments

VDJ	Either (for platypus version "v2") output from VDJ_analyze function. This should be a list of clonotype dataframes, with each list element corresponding to a single VDJ repertoire, OR (for platypus version "v3") the the VDJ matrix output of the VDJ_GEX_matrix() function (VDJ.GEX.matrix.output[[1]])
HC.gene.number	Numeric value indicating the top genes to be displayed. If this number is higher than the total number of unique HC V genes in the VDJ repertoire, then this number is equal to the number of unique HC V genes.
LC.Vgene	Logical indicating whether to make a barplot of the LC V genes distribution. Default is set to FALSE.
LC.gene.number	Numeric value indicating the top genes to be displayed. If this number is higher than the total number of unique LC V genes in the VDJ repertoire, then this number is equal to the number of unique LC V genes.
platypus.version	Character. Defaults to "v3". Can be "v2" or "v3" dependent on the input format

Value

Returns a list of ggplot objects which show the distribution of IgH and IgK/L V genes for the most used V genes.

Examples

```
VDJ_Vgene_usage_barplot(VDJ = Platypus::small_vgm[[1]],
  HC.gene.number = 2, platypus.version = "v3")
```

VDJ_Vgene_usage_stacked_barplot

Produces a stacked barplot with the fraction of the most frequently used IgH and IgK/L Vgenes. This function can be used in combination with the VDJ_Vgene_usage_barplot to visualize V gene usage per sample and among samples.

Description

Produces a stacked barplot with the fraction of the most frequently used IgH and IgK/L Vgenes. This function can be used in combination with the VDJ_Vgene_usage_barplot to visualize V gene usage per sample and among samples.

Usage

```
VDJ_Vgene_usage_stacked_barplot(
  VDJ,
  HC.gene.number,
  Fraction.HC,
  LC.Vgene,
  LC.gene.number,
  Fraction.LC,
  platypus.version
)
```

Arguments

VDJ	Either (for platypus version "v2") output from VDJ_analyze function. This should be a list of clonotype dataframes, with each list element corresponding to a single VDJ repertoire, OR (for platypus version "v3") the the VDJ matrix output of the VDJ_GEX_matrix() function (normally VDJ.GEX.matrix.output[[1]])
HC.gene.number	Numeric value indicating the top genes to be displayed. If this number is higher than the total number of unique HC V genes in the VDJ repertoire, then this number is equal to the number of unique HC V genes.
Fraction.HC	Numeric value indicating the minimum fraction of clones expressing a particular HC V gene. If the usage of a particular gene is below this value, then this gene is excluded. If the usage of a particular gene is above this value even in one sample, then this gene is included in the analysis. Default value is set to 0, thus all genes are selected.
LC.Vgene	Logical indicating whether to make a barplot of the LC V gene distribution. Default is set to FALSE.

LC.gene.number Numeric value indicating the top genes to be displayed. If this number is higher than the total number of unique LC V genes in the VDJ repertoire, then this number is equal to the number of unique LC V genes.

Fraction.LC Numeric value indicating the minimum fraction of clones expressing a particular LC V gene. If the usage of a particular gene is below this value, then this gene is excluded. If the usage of a particular gene is above this value even in one sample, then this gene is included in the analysis. Default value is set to 0, thus all genes are selected.

platypus.version Set according to input format to either "v2" or "v3". Defaults to "v3"

Value

Returns a list of ggplot objects which show the stacked distribution of IgH and IgK/L V genes for the most used V genes. Returns an empty plot if the Fraction.HC or Fraction.LC that were selected were too high, resulting in the exclusion of all the genes.

Examples

```
#Platypus v3
example.vdj.vgene_usage <- VDJ_Vgene_usage_stacked_barplot(
VDJ = Platypus::small_vgm[[1]], LC.Vgene = TRUE
,HC.gene.number = 15, Fraction.HC = 1, platypus.version = "v3")
```

VDJ_VJ_usage_circos *Makes a Circos plot from the VDJ_analyze output. Connects the V gene with the corresponding J gene for each clonotype.*

Description

Makes a Circos plot from the VDJ_analyze output. Connects the V gene with the corresponding J gene for each clonotype.

Usage

```
VDJ_VJ_usage_circos(
  VDJ,
  A.or.B,
  label.threshold,
  cell.level,
  c.threshold,
  clonotype.per.gene.threshold,
  c.count,
  platypus.version,
  filter1H1L
)
```

Arguments

VDJ	The output of the VDJ_GEX_integrate function (Platypus platypus.version v2). A list of data frames for each sample containing the clonotype information and cluster membership information. For Platypus platypus.version v3, the VDJ output of the VDJ_GEX_matrix function (VDJ_GEX_matrix.output[[1]]) has to be supplied.
A.or.B	Determines whether to plot the V J gene pairing of the alpha or beta chain. "A", "B" or "both" as possible inputs. Default: "both".
label.threshold	Minimal amount of clonotypes per gene necessary to add a gene label to the sector. Default: 0.
cell.level	Logical, defines whether weight of connection should be based on number of clonotypes or number of cells. Default: number of clonotypes.
c.threshold	Only clonotypes are considered with a frequency higher then c.threshold. Allows to filter for only highly expanded clonotypes.
clonotype.per.gene.threshold	How many clonotypes are required to plot a sector for a gene. Filters the rows and columns of the final adjacency matrix.
c.count	Show clonotype or cell count on Circos plot. Default = T.
platypus.version	Which platypus.version of platypus is being used. Default = v3. Set to v3 if VDJ_GEX_matrix.output[[1]] is used
filter1H1L	Whether to filter the input VDJ in "v3" to only include cells with 1 VDJ and 1 VJ chain. Defaults to TRUE

Value

Returns list of plots. The first n elements contain the circos plot of the n datasets from the VDJ.analyze function. The n+1 element contains a list of the n adjacency matrices for each dataset.

Examples

```
plots <- VDJ_VJ_usage_circos(VDJ = Platypus::small_vgm[[1]], platypus.version = "v3",
cell.level = TRUE)
```

Index

* datasets

Bcell_sequences_example_tree, 24
Bcell_tree_2, 25
class_switch_prob_hum, 26
class_switch_prob_mus, 27
colors, 31
hotspot_df, 64
hum_b_h, 65
hum_b_l, 65
hum_t_h, 66
hum_t_l, 67
iso_SHM_prob, 67
mus_b_h, 68
mus_b_l, 68
mus_b_trans, 69
mus_t_h, 70
mus_t_l, 71
new, 71
one_spot_df, 72
pheno_SHM_prob, 73
small_vgm, 87
special_v, 88
trans_switch_prob_b, 88
trans_switch_prob_t, 89
vdj_length_prob, 125

AbForests_AntibodyForest, 4
AbForests_CompareForests, 7
AbForests_ConvertStructure, 10
AbForests_CsvToDf, 11
AbForests_ForestMetrics, 11
AbForests_PlotGraphs, 14
AbForests_PlyloToMatrix, 15
AbForests_RemoveNets, 16
AbForests_SubRepertoiresByCells, 18
AbForests_SubRepertoiresByUniqueSeq,
19
AbForests_UniqueAntibodyVariants, 21
automate_GEX, 22

Bcell_sequences_example_tree, 24
Bcell_tree_2, 25

call_MIXCR, 25
class_switch_prob_hum, 26
class_switch_prob_mus, 27
clonofreq, 27
clonofreq.isotype.data, 28
clonofreq.isotype.plot, 28
clonofreq.trans.data, 29
clonofreq.trans.plot, 30
cluster.id.igraph, 30
colors, 31

Echidna_vae_generate, 32

get.avr.mut.data, 33
get.avr.mut.plot, 33
get.barplot.errorbar, 34
get.elbow, 35
get.n.node.data, 35
get.n.node.plot, 36
get.seq.distance, 36
get.umap, 37
get.vgu.matrix, 37
GEX_automate, 38
GEX_clonotype, 40
GEX_cluster_genes, 41
GEX_cluster_genes_heatmap, 42
GEX_cluster_membership, 43
GEX_coexpression_coefficient, 44
GEX_DEgenes, 45
GEX_DEgenes_persample, 48
GEX_dottile_plot, 49
GEX_G0term, 50
GEX_GSEA, 52
GEX_heatmap, 54
GEX_pairwise_DEGs, 55
GEX_phenotype, 56
GEX_phenotype_per_clone, 57

GEX_proportions_barplot, 58
GEX_scatter_coexpression, 59
GEX_topN_DE_genes_per_cluster, 60
GEX_visualize_clones, 61
GEX_volcano, 62

hotspot_df, 64
hum_b_h, 65
hum_b_l, 65
hum_t_h, 66
hum_t_l, 67

iso_SHM_prob, 67

mus_b_h, 68
mus_b_l, 68
mus_b_trans, 69
mus_t_h, 70
mus_t_l, 71

new, 71
no.empty.node, 72

one_spot_df, 72

pheno_SHM_prob, 73
PlatypusDB_AIRR_to_VGM, 74
PlatypusDB_fetch, 75
PlatypusDB_find_CDR3s, 77
PlatypusDB_list_projects, 78
PlatypusDB_load_from_disk, 79
PlatypusDB_VGM_to_AIRR, 80

select.top, 82
simulate_repertoire, 83
small_vgm, 87
special_v, 88

trans_switch_prob_b, 88
trans_switch_prob_t, 89

umap.top.highlight, 89

VDJ_alpha_beta_Vgene_circos, 90
VDJ_analyze, 91
VDJ_assemble_for_PnP, 93
VDJ_call_MIXCR, 95
VDJ_circos, 96
VDJ_clonal_donut, 97
VDJ_clonal_expansion, 99
VDJ_clonal_lineages, 101
VDJ_clonotype, 102
VDJ_clonotype_clusters_circos, 104
VDJ_contigs_to_vgm, 106
VDJ_diversity, 106
VDJ_dublets, 108
VDJ_extract_germline, 109
VDJ_GEX_clonal_lineage_clusters, 110
VDJ_GEX_expansion, 111
VDJ_GEX_integrate, 112
VDJ_GEX_matrix, 113
VDJ_GEX_overlay_clones, 119
VDJ_GEX_stats, 122
VDJ_isotypes_per_clone, 123
vdj_length_prob, 124
VDJ_logoplot_vector, 126
VDJ_network, 126
VDJ_overlap_heatmap, 128
VDJ_per_clone, 129
VDJ_plot_SHM, 130
VDJ_reclonotype_list_arrange, 132
VDJ_tree, 133
VDJ_variants_per_clone, 134
VDJ_Vgene_usage, 135
VDJ_Vgene_usage_barplot, 136
VDJ_Vgene_usage_stacked_barplot, 137
VDJ_VJ_usage_circos, 138