

# Package: coMMpass (via r-universe)

May 13, 2026

**Title** MMRF CoMMpass Data Analysis Pipeline

**Version** 0.1.0

**Description** A reproducible data acquisition and analysis pipeline for the MMRF CoMMpass study using Targets and Nix.

**License** MIT + file LICENSE

**URL** <https://JohnGavin.github.io/coMMpass-analysis/>,  
<https://github.com/JohnGavin/coMMpass-analysis>

**BugReports** <https://github.com/JohnGavin/coMMpass-analysis/issues>

**Encoding** UTF-8

**RoxygenNote** 7.3.3

**VignetteBuilder** quarto, knitr

**Imports** TCGAbiolinks, S4Vectors, SummarizedExperiment, edgeR, logger, dplyr, aws.s3, arrow, DBI, duckdb, dbplyr, cli, jsonlite, glue, stringr, tibble

**Suggests** anndataR, AnnotationDbi, SingleCellExperiment, apeglm, quarto, targets, tarchetypes, crew, DESeq2, fgsea, GenomicDataCommons, ggplot2, limma, msigdb, org.Hs.eg.db, scales, survival, tidyverse, knitr, rmarkdown, pkgdown, testthat (>= 3.0.0), gert, gh, nanoparquet, DT, htmltools, plotly, plumber, visNetwork, here, dagitty, ggdag, naniar, brms, bayesplot, posterior, splines2

**biocViews** Software

**Additional\_repositories** <https://bioc.r-universe.dev>

**Config/testthat/edition** 3

**Config/pak/sysreqs**

cmake libicu-dev libpng-dev libxml2-dev libssl-dev libx11-dev xz-utils zlib1g-dev

**Repository** <https://johngavin.r-universe.dev>

**Date/Publication** 2026-05-13 16:51:00 UTC

**RemoteUrl** <https://github.com/JohnGavin/coMMpass-analysis>

**RemoteRef** HEAD

**RemoteSha** 78011582e659df92455d7cb1c50c5c329c359f3c

## Contents

acquire_compass_data . . . . .	4
annotate_de_results . . . . .	5
annotate_genes . . . . .	6
api_get_clinical . . . . .	6
api_get_de_results . . . . .	7
api_get_pathways . . . . .	8
api_get_survival . . . . .	8
api_list_datasets . . . . .	9
api_serve . . . . .	9
calculate_cooccurrence . . . . .	10
calculate_qc_metrics . . . . .	11
check_adjustment . . . . .	11
check_dependencies . . . . .	12
clean_clinical_data . . . . .	12
clean_expression_data . . . . .	13
clean_treatment_data . . . . .	13
compass_dag . . . . .	14
compute_pca . . . . .	14
compute_riss . . . . .	15
correlate_genes . . . . .	16
correlate_genes_batch . . . . .	16
create_project_dirs . . . . .	17
create_summary_table . . . . .	17
download_aws_data . . . . .	18
download_clinical_data . . . . .	19
download_gdc_rnaseq . . . . .	20
download_s3_subset . . . . .	21
example_data . . . . .	22
export_h5ad . . . . .	23
extract_cytogenetic_data . . . . .	24
extract_risk_table . . . . .	24
filter_low_quality . . . . .	25
find_consensus_genes . . . . .	26
format_file_size . . . . .	26
format_with_commas . . . . .	27
gene_report . . . . .	27
generate_api_endpoint . . . . .	28
generate_api_index . . . . .	29
generate_summary_report . . . . .	30
get_adjustment_sets . . . . .	31
get_compass_clinical . . . . .	31
get_compass_data_dictionary . . . . .	32
get_compass_tbl . . . . .	33
get_counts_assay . . . . .	34
get_variable_docs . . . . .	34
integrate_clinical_expression . . . . .	35

list_s3_compass . . . . .	35
normalize_rnaseq . . . . .	36
plot_cooccurrence_heatmap . . . . .	37
plot_cytogenetic_oncoprint . . . . .	37
plot_dag . . . . .	38
plot_enrichment_barplot . . . . .	39
plot_enrichment_dotplot . . . . .	40
plot_expression_by_subtype . . . . .	40
plot_forest . . . . .	41
plot_gene_correlation . . . . .	42
plot_gsea_running_score . . . . .	42
plot_heatmap_de . . . . .	43
plot_km . . . . .	44
plot_ma . . . . .	45
plot_pca . . . . .	45
plot_volcano . . . . .	46
prepare_survival_data . . . . .	47
query_compass_parquet . . . . .	47
query_compass_rna . . . . .	49
render_de_report . . . . .	49
run_cox_regression . . . . .	50
run_deseq2 . . . . .	51
run_edger . . . . .	52
run_gsea . . . . .	52
run_kaplan_meier . . . . .	53
run_km_by_expression . . . . .	54
run_km_by_markers . . . . .	55
run_limma . . . . .	55
run_ora . . . . .	56
run_pathway_analysis . . . . .	57
run_vst . . . . .	57
save_timestamped . . . . .	58
setup_logging . . . . .	58
strip_plotly . . . . .	59
summarize_cytogenetics . . . . .	59
summarize_data . . . . .	60
summarize_de_methods . . . . .	61
summarize_treatment . . . . .	61

---

acquire\_commpass\_data *Main data acquisition function*

---

### Description

Orchestrates the download of CoMMpass data from various sources including RNA-seq data from GDC, clinical data, and optionally AWS data.

### Usage

```
acquire_commpass_data(  
  download_rnaseq = TRUE,  
  download_clinical = TRUE,  
  download_aws = FALSE,  
  sample_limit = 200,  
  random_sample = TRUE,  
  seed = 42,  
  use_parquet = TRUE  
)
```

### Arguments

download_rnaseq	Whether to download RNA-seq data
download_clinical	Whether to download clinical data
download_aws	Whether to download from AWS
sample_limit	Limit number of samples (NULL for all)
random_sample	If TRUE, randomly sample patients
seed	Random seed for sampling
use_parquet	If TRUE, save parquet files

### Value

List of file paths to the downloaded data

### See Also

Other data-acquisition: [download\\_clinical\\_data\(\)](#), [download\\_gdc\\_rnaseq\(\)](#), [download\\_s3\\_subset\(\)](#), [get\\_commpass\\_clinical\(\)](#), [list\\_s3\\_commpass\(\)](#), [query\\_commpass\\_rna\(\)](#)

## Examples

```
## Not run:  
# Download only clinical data  
results <- acquire_compass_data(  
  download_rnaseq = FALSE,  
  download_clinical = TRUE,  
  download_aws = FALSE  
)  
  
## End(Not run)
```

---

annotate\_de\_results    *Add gene symbol column to DE results table*

---

## Description

Add gene symbol column to DE results table

## Usage

```
annotate_de_results(de_results, annotation)
```

## Arguments

de_results	Data frame with Ensembl IDs as rownames
annotation	Data frame from <code>annotate_genes()</code>

## Value

DE results with `gene_symbol` column prepended

## See Also

Other pathway: [annotate\\_genes\(\)](#), [plot\\_enrichment\\_barplot\(\)](#), [plot\\_enrichment\\_dotplot\(\)](#), [plot\\_gsea\\_running\\_score\(\)](#), [run\\_gsea\(\)](#), [run\\_ora\(\)](#), [run\\_pathway\\_analysis\(\)](#)

---

annotate\_genes      *Annotate Ensembl gene IDs with symbols and descriptions*

---

### Description

Maps Ensembl gene IDs to HGNC symbols and Entrez IDs. Uses msigdb gene mapping as the primary source (always available with msigdb). Falls back to org.Hs.eg.db + AnnotationDbi if msigdb is not installed.

### Usage

```
annotate_genes(ensembl_ids)
```

### Arguments

ensembl\_ids      Character vector of Ensembl gene IDs (versions stripped)

### Value

Data frame with columns: ensembl\_gene, gene\_symbol, entrez\_gene. Unmappable IDs have NA for symbol/entrez.

### See Also

Other pathway: [annotate\\_de\\_results\(\)](#), [plot\\_enrichment\\_barplot\(\)](#), [plot\\_enrichment\\_dotplot\(\)](#), [plot\\_gsea\\_running\\_score\(\)](#), [run\\_gsea\(\)](#), [run\\_ora\(\)](#), [run\\_pathway\\_analysis\(\)](#)

---

api\_get\_clinical      *Get clinical data for API response*

---

### Description

Reads clinical data from the targets store and formats it for API delivery. Optionally filters by patient IDs or variables.

### Usage

```
api_get_clinical(patient_ids = NULL, variables = NULL, store = "_targets")
```

### Arguments

patient\_ids      Optional character vector of patient IDs to filter  
 variables      Optional character vector of column names to select  
 store          Path to targets store (default: "\_targets")

**Value**

A data frame of clinical data, or NULL if unavailable

**See Also**

Other api: [api\\_get\\_de\\_results\(\)](#), [api\\_get\\_pathways\(\)](#), [api\\_get\\_survival\(\)](#), [api\\_list\\_datasets\(\)](#), [api\\_serve\(\)](#), [generate\\_api\\_endpoint\(\)](#), [generate\\_api\\_index\(\)](#)

---

api_get_de_results	<i>Get DE results for API response</i>
--------------------	--

---

**Description**

Reads differential expression results from the targets store. Returns the DESeq2 results table with gene symbols.

**Usage**

```
api_get_de_results(padj_threshold = 1, lfc_threshold = 0, store = "_targets")
```

**Arguments**

padj_threshold	Filter to genes with padj below this threshold (default: 1, i.e., no filter)
lfc_threshold	Filter to genes with absolute log2FC above this threshold (default: 0, i.e., no filter)
store	Path to targets store

**Value**

A data frame of DE results, or NULL if unavailable

**See Also**

Other api: [api\\_get\\_clinical\(\)](#), [api\\_get\\_pathways\(\)](#), [api\\_get\\_survival\(\)](#), [api\\_list\\_datasets\(\)](#), [api\\_serve\(\)](#), [generate\\_api\\_endpoint\(\)](#), [generate\\_api\\_index\(\)](#)

---

api_get_pathways	<i>Get pathway analysis results for API response</i>
------------------	--

---

**Description**

Reads GSEA results from the targets store.

**Usage**

```
api_get_pathways(significant_only = FALSE, store = "_targets")
```

**Arguments**

significant_only	If TRUE, return only pathways with padj < 0.05
store	Path to targets store

**Value**

A data frame of pathway results, or NULL if unavailable

**See Also**

Other api: [api\\_get\\_clinical\(\)](#), [api\\_get\\_de\\_results\(\)](#), [api\\_get\\_survival\(\)](#), [api\\_list\\_datasets\(\)](#), [api\\_serve\(\)](#), [generate\\_api\\_endpoint\(\)](#), [generate\\_api\\_index\(\)](#)

---

api_get_survival	<i>Get survival data for API response</i>
------------------	---

---

**Description**

Reads prepared survival data from the targets store.

**Usage**

```
api_get_survival(store = "_targets")
```

**Arguments**

store	Path to targets store
-------	-----------------------

**Value**

A data frame of survival data, or NULL if unavailable

**See Also**

Other api: [api\\_get\\_clinical\(\)](#), [api\\_get\\_de\\_results\(\)](#), [api\\_get\\_pathways\(\)](#), [api\\_list\\_datasets\(\)](#), [api\\_serve\(\)](#), [generate\\_api\\_endpoint\(\)](#), [generate\\_api\\_index\(\)](#)

---

api\_list\_datasets      *List available API datasets*

---

**Description**

Returns metadata about all datasets available through the API, including data type, row/column counts, and download formats.

**Usage**

```
api_list_datasets()
```

**Value**

A data frame with columns: dataset, description, format, n\_rows, n\_cols, last\_updated

**See Also**

Other api: [api\\_get\\_clinical\(\)](#), [api\\_get\\_de\\_results\(\)](#), [api\\_get\\_pathways\(\)](#), [api\\_get\\_survival\(\)](#), [api\\_serve\(\)](#), [generate\\_api\\_endpoint\(\)](#), [generate\\_api\\_index\(\)](#)

**Examples**

```
## Not run:  
api_list_datasets()  
  
## End(Not run)
```

---

api\_serve      *Launch the plumber API*

---

**Description**

Starts the plumber API server for programmatic data access. The API reads pre-computed results from the targets store.

**Usage**

```
api_serve(port = 8080L, host = "127.0.0.1", store = "_targets")
```

**Arguments**

port	Port number (default: 8080)
host	Host address (default: "127.0.0.1")
store	Path to targets store (default: "_targets")

**Value**

Invisible NULL (starts server)

**See Also**

Other api: [api\\_get\\_clinical\(\)](#), [api\\_get\\_de\\_results\(\)](#), [api\\_get\\_pathways\(\)](#), [api\\_get\\_survival\(\)](#), [api\\_list\\_datasets\(\)](#), [generate\\_api\\_endpoint\(\)](#), [generate\\_api\\_index\(\)](#)

**Examples**

```
## Not run:
# Start the API server
api_serve(port = 8080)

# Then query from R or curl:
# curl http://localhost:8080/datasets
# curl http://localhost:8080/data/clinical?variables=submitter_id,gender

## End(Not run)
```

---

calculate\_cooccurrence

*Calculate pairwise co-occurrence of cytogenetic alterations*

---

**Description**

Computes Fisher's exact test for each pair of cytogenetic markers to identify significant co-occurrence or mutual exclusivity patterns.

**Usage**

```
calculate_cooccurrence(cyto_data, markers = NULL)
```

**Arguments**

cyto_data	Data frame from <a href="#">extract_cytogenetic_data()</a>
markers	Character vector of marker columns. Default auto-detects.

**Value**

Data frame with columns: marker1, marker2, odds\_ratio, pvalue, padj (BH-corrected), n\_both, n\_either, tendency (co-occurrence/exclusive/none)

**See Also**

Other cytogenetics: [compute\\_riss\(\)](#), [extract\\_cytogenetic\\_data\(\)](#), [plot\\_cooccurrence\\_heatmap\(\)](#), [plot\\_cytogenetic\\_oncoprint\(\)](#), [plot\\_expression\\_by\\_subtype\(\)](#), [summarize\\_cytogenetics\(\)](#)

**Examples**

```
## Not run:
cyto <- arrow::read_parquet("data/raw/clinical/cytogenetic_data.parquet")
cooc <- calculate_cooccurrence(cyto)

## End(Not run)
```

---

calculate\_qc\_metrics    *Calculate QC metrics for RNA-seq data*

---

**Description**

Calculate QC metrics for RNA-seq data

**Usage**

```
calculate_qc_metrics(se_data)
```

**Arguments**

se\_data            A SummarizedExperiment object with a "counts" assay

**Value**

A data frame with QC metrics per sample

---

check\_adjustment    *Compare model covariates against DAG-implied adjustments*

---

**Description**

Checks whether the covariates used in a model match any of the DAG-implied minimal sufficient adjustment sets.

**Usage**

```
check_adjustment(
  model_covariates,
  exposure = "cytogenetic_risk",
  outcome = "overall_survival",
  dag = NULL
)
```

**Arguments**

model_covariates	Character vector of covariates used in the model.
exposure	Character. Exposure variable.
outcome	Character. Outcome variable.
dag	A ‘dagitty’ DAG object.

**Value**

A list with ‘\$sufficient’ (logical), ‘\$model\_covariates’, ‘\$adjustment\_sets’, and ‘\$recommendation’.

---

check_dependencies	<i>Check package dependencies</i>
--------------------	-----------------------------------

---

**Description**

Check package dependencies

**Usage**

```
check_dependencies()
```

---

clean_clinical_data	<i>Clean Clinical Data</i>
---------------------	----------------------------

---

**Description**

Standardizes clinical data column names and formats

**Usage**

```
clean_clinical_data(clinical_raw)
```

**Arguments**

clinical_raw	Raw clinical data frame
--------------	-------------------------

**Value**

Cleaned clinical data frame

**See Also**

Other data-cleaning: [clean\\_expression\\_data\(\)](#), [clean\\_treatment\\_data\(\)](#), [integrate\\_clinical\\_expression\(\)](#), [summarize\\_treatment\(\)](#)

---

clean\_expression\_data *Clean Expression Data*

---

**Description**

Standardizes expression data format and adds metadata

**Usage**

```
clean_expression_data(expr_raw)
```

**Arguments**

expr\_raw            Raw expression data (matrix or data frame)

**Value**

Cleaned expression matrix with gene names as rownames

**See Also**

Other data-cleaning: [clean\\_clinical\\_data\(\)](#), [clean\\_treatment\\_data\(\)](#), [integrate\\_clinical\\_expression\(\)](#), [summarize\\_treatment\(\)](#)

---

clean\_treatment\_data *Clean and standardize treatment data*

---

**Description**

Standardizes regimen names, derives regimen class, and creates an ordered response factor from raw CoMMpass treatment response data.

**Usage**

```
clean_treatment_data(trtresp_raw)
```

**Arguments**

trtresp\_raw        Data frame of raw treatment response data with columns including patient/submitter ID, treatment line, regimen description, best response, and transplant status

**Value**

Cleaned data frame with standardized columns: 'patient\_id', 'treatment\_line', 'regimen\_name', 'regimen\_class', 'best\_response' (ordered factor), 'stem\_cell\_transplant' (logical), 'treatment\_start\_days' (integer)

**See Also**

Other data-cleaning: [clean\\_clinical\\_data\(\)](#), [clean\\_expression\\_data\(\)](#), [integrate\\_clinical\\_expression\(\)](#), [summarize\\_treatment\(\)](#)

---

compass_dag	<i>Define the CoMMpass causal DAG</i>
-------------	---------------------------------------

---

**Description**

Encodes domain-knowledge causal assumptions for the CoMMpass multiple myeloma analysis. Nodes represent measured or latent variables; edges represent assumed causal effects.

**Usage**

```
compass_dag()
```

**Value**

A ‘dagitty::dagitty’ DAG object

---

compute_pca	<i>Compute PCA from transformed expression data</i>
-------------	---

---

**Description**

Compute PCA from transformed expression data

**Usage**

```
compute_pca(expr_matrix, n_top = 500L, metadata = NULL)
```

**Arguments**

expr_matrix	Numeric matrix (genes x samples) of transformed expression values (e.g., VST or logCPM)
n_top	Number of most variable genes to use (default 500)
metadata	Optional data frame of sample annotations to merge with PCA coordinates

**Value**

List with components: coords (data frame with PC1, PC2, ...), var\_explained (numeric vector of variance proportions), n\_genes\_used

**See Also**

Other differential-expression: [plot\\_heatmap\\_de\(\)](#), [plot\\_ma\(\)](#), [plot\\_pca\(\)](#), [plot\\_volcano\(\)](#), [run\\_deseq2\(\)](#), [run\\_vst\(\)](#), [summarize\\_de\\_methods\(\)](#)

---

compute_riss	<i>Compute Revised International Staging System (R-ISS)</i>
--------------	---

---

### Description

Classifies patients into R-ISS stages based on ISS stage, cytogenetic risk group, and LDH level (Palumbo et al., JCO 2015).

### Usage

```
compute_riss(iss_stage, risk_group, ldh = NA_real_, ldh_uln = 250)
```

### Arguments

iss_stage	Character vector: "Stage I", "Stage II", "Stage III"
risk_group	Character vector: "High", "Standard", "high", "standard"
ldh	Numeric vector: LDH value (U/L), or NA
ldh_uln	Numeric scalar: upper limit of normal for LDH (default 250)

### Value

Character vector: "R-ISS I", "R-ISS II", "R-ISS III", or NA

### See Also

Other cytogenetics: [calculate\\_cooccurrence\(\)](#), [extract\\_cytogenetic\\_data\(\)](#), [plot\\_cooccurrence\\_heatmap\(\)](#), [plot\\_cytogenetic\\_oncoprint\(\)](#), [plot\\_expression\\_by\\_subtype\(\)](#), [summarize\\_cytogenetics\(\)](#)

### Examples

```
compute_riss("Stage I", "Standard", ldh = 200)
compute_riss("Stage III", "High", ldh = 300)
compute_riss(
  c("Stage I", "Stage III", "Stage II"),
  c("Standard", "High", "Standard"),
  ldh = c(200, 300, NA)
)
```

---

correlate\_genes      *Correlate two genes across samples*

---

### Description

Computes Pearson or Spearman correlation between two genes in an expression matrix.

### Usage

```
correlate_genes(expr_matrix, gene_x, gene_y, method = c("pearson", "spearman"))
```

### Arguments

expr_matrix	Numeric matrix (genes x samples) of transformed expression values (e.g., VST)
gene_x	Gene name for x-axis
gene_y	Gene name for y-axis
method	Correlation method: "pearson" (default) or "spearman"

### Value

List with components: gene\_x, gene\_y, method, estimate, p\_value, n\_samples, expr\_x (numeric vector), expr\_y (numeric vector)

### See Also

Other gene-correlation: [correlate\\_genes\\_batch\(\)](#), [plot\\_gene\\_correlation\(\)](#)

---

correlate\_genes\_batch      *Batch correlation: one gene vs many*

---

### Description

Correlates a target gene against a vector of candidate genes and returns a summary data frame sorted by absolute correlation.

### Usage

```
correlate_genes_batch(
  expr_matrix,
  target_gene,
  candidate_genes,
  method = c("pearson", "spearman")
)
```

**Arguments**

expr_matrix	Numeric matrix (genes x samples)
target_gene	Gene name to correlate against
candidate_genes	Character vector of gene names to test
method	Correlation method: "pearson" or "spearman"

**Value**

Data frame with columns: gene, estimate, p\_value, padj, n\_samples, method

**See Also**

Other gene-correlation: [correlate\\_genes\(\)](#), [plot\\_gene\\_correlation\(\)](#)

---

create\_project\_dirs    *Create project directories*

---

**Description**

Create project directories

**Usage**

```
create_project_dirs(base_dir = ".")
```

**Arguments**

base_dir	Base directory path (default: ".")
----------	------------------------------------

---

create\_summary\_table    *Create Summary Statistics Table*

---

**Description**

Generate a nicely formatted summary statistics table for numeric variables

**Usage**

```
create_summary_table(data, vars = NULL)
```

**Arguments**

data	Data frame
vars	Character vector of variable names (NULL for all numeric)

**Value**

Data frame with summary statistics

**See Also**

Other utilities: [example\\_data\(\)](#), [export\\_h5ad\(\)](#), [format\\_file\\_size\(\)](#), [format\\_with\\_commas\(\)](#), [gene\\_report\(\)](#), [strip\\_plotly\(\)](#)

---

download_aws_data	<i>Download data from AWS S3 open access bucket</i>
-------------------	---

---

**Description**

Download data from AWS S3 open access bucket

**Usage**

```
download_aws_data(  
  bucket_name = "gdc-mmrf-commpass-phs000748-2-open",  
  prefix = NULL,  
  data_dir = "data/raw/aws",  
  region = "us-east-1"  
)
```

**Arguments**

bucket_name	S3 bucket name
prefix	File prefix to filter
data_dir	Local directory for downloads
region	AWS region

**Value**

Directory path as character string

---

`download_clinical_data`*Download clinical data from GDC*

---

**Description**

Downloads clinical and biospecimen data from the Genomic Data Commons (GDC) for the specified project. Data is saved in parquet and RDS formats.

**Usage**

```
download_clinical_data(  
  project_id = "MMRF-COMMPASS",  
  data_dir = "data/raw/clinical",  
  use_parquet = TRUE  
)
```

**Arguments**

<code>project_id</code>	Project identifier (default: "MMRF-COMMPASS")
<code>data_dir</code>	Directory to save data
<code>use_parquet</code>	If TRUE, also save parquet files (default TRUE)

**Value**

Path to the directory containing the saved data files

**See Also**

Other data-acquisition: [acquire\\_commpass\\_data\(\)](#), [download\\_gdc\\_rnaseq\(\)](#), [download\\_s3\\_subset\(\)](#), [get\\_commpass\\_clinical\(\)](#), [list\\_s3\\_commpass\(\)](#), [query\\_commpass\\_rna\(\)](#)

**Examples**

```
## Not run:  
# Download clinical data  
clinical_dir <- download_clinical_data()  
  
# Load from parquet  
clinical <- arrow::read_parquet(file.path(clinical_dir, "clinical_data.parquet"))  
  
## End(Not run)
```

---

download\_gdc\_rnaseq     *Download RNA-seq data from GDC*

---

## Description

Downloads RNA-seq gene expression data from the Genomic Data Commons (GDC) for the specified project. Data is saved as a SummarizedExperiment RDS and as parquet files (counts, sample metadata, gene metadata).

## Usage

```
download_gdc_rnaseq(  
  project_id = "MMRF-COMMPASS",  
  data_dir = "data/raw/gdc",  
  sample_limit = 200,  
  random_sample = TRUE,  
  seed = 42,  
  use_parquet = TRUE  
)
```

## Arguments

project_id	Project identifier (default: "MMRF-COMMPASS")
data_dir	Directory to save data
sample_limit	Maximum number of samples (NULL for all, default 200)
random_sample	If TRUE and sample_limit is set, randomly sample patients using seed for reproducibility (default TRUE)
seed	Random seed for reproducible sampling (default 42)
use_parquet	If TRUE, also save parquet files alongside RDS (default TRUE)

## Value

Path to the saved RDS file containing the SummarizedExperiment

## See Also

Other data-acquisition: [acquire\\_commpass\\_data\(\)](#), [download\\_clinical\\_data\(\)](#), [download\\_s3\\_subset\(\)](#), [get\\_commpass\\_clinical\(\)](#), [list\\_s3\\_commpass\(\)](#), [query\\_commpass\\_rna\(\)](#)

## Examples

```
## Not run:  
# Download 200 random samples with parquet output  
rnaseq_file <- download_gdc_rnaseq(sample_limit = 200)  
  
# Load the SummarizedExperiment
```

```
se_data <- readRDS(rnaseq_file)

# Or load parquet directly
counts <- arrow::read_parquet("data/raw/gdc/rnaseq_counts.parquet")

## End(Not run)
```

---

download\_s3\_subset      *Download a Sample of RNA-seq Files from S3*

---

## Description

Downloads a subset of RNA-seq files from the public MMRF CoMMpass S3 bucket. This function is useful for testing and development with a small sample of data before downloading the full dataset. Files are downloaded using anonymous access to the public bucket.

## Usage

```
download_s3_subset(s3_paths, dest_dir = "data/raw/rna_seq", n = 3)
```

## Arguments

s3_paths	Character vector of S3 object keys (file paths) to download from
dest_dir	Destination directory for downloaded files (default: "data/raw/rna_seq")
n	Number of files to download (default: 3)

## Value

Character vector of successfully downloaded file paths

## See Also

Other data-acquisition: [acquire\\_commpass\\_data\(\)](#), [download\\_clinical\\_data\(\)](#), [download\\_gdc\\_rnaseq\(\)](#), [get\\_commpass\\_clinical\(\)](#), [list\\_s3\\_commpass\(\)](#), [query\\_commpass\\_rna\(\)](#)

## Examples

```
## Not run:
# List available files
s3_files <- list_s3_commpass()

# Download first 3 RNA-seq files
downloaded <- download_s3_subset(s3_files, n = 3)

# Check what was downloaded
basename(downloaded)

## End(Not run)
```

---

 example\_data

*Load example coMMpass datasets*


---

## Description

Returns small **synthetic** datasets (no real patient data) for interactive exploration and testing. The data exercises the full pipeline: QC, cleaning, survival analysis, and cytogenetic classification.

## Usage

```
example_data()
```

## Details

The SummarizedExperiment is reconstructed at load time from stored plain-R components (matrix + data.frames), so the RDS files have no S4 class dependency.

## Value

A named list with three elements:

**rnaseq\_se** A [SummarizedExperiment::SummarizedExperiment] with 50 genes x 20 samples. Assay named "unstranded" (matches GDC format). rowData contains 'gene\_id' (Ensembl-format with version suffix) and 'gene\_type'. colData contains 'submitter\_id' and 'sample\_type'.

**clinical** A data.frame (20 patients) with columns 'submitter\_id', 'vital\_status', 'days\_to\_death', 'days\_to\_last\_follow\_up', 'age\_at\_diagnosis' (in days), 'gender', 'iss\_stage', 'heavy\_chain' (myeloma isotype), 'light\_chain' (Kappa/Lambda), 'ecog\_status' (0-4), 'ldh' (U/L), 'b2m' (beta-2 microglobulin, mg/L), 'albumin' (g/dL), 'flc\_kappa' (free kappa light chain, mg/L), 'flc\_lambda' (free lambda light chain, mg/L), 'hemoglobin' (g/dL), 'creatinine' (mg/dL), 'calcium' (corrected, mg/dL), 'platelets' (10<sup>9</sup>/L).

**cytogenetic** A data.frame (20 patients) with columns 'patient\_id', 't\_4\_14', 't\_11\_14', 't\_14\_16', 'del\_17p', 'gain\_1q', 'risk\_group'.

**treatment** A data.frame of treatment lines (1-3 per patient) with columns 'patient\_id', 'treatment\_line', 'regimen\_name', 'regimen\_class', 'best\_response' (ordered factor), 'stem\_cell\_transplant' (logical), 'treatment\_start\_days' (integer).

## See Also

Other utilities: [create\\_summary\\_table\(\)](#), [export\\_h5ad\(\)](#), [format\\_file\\_size\(\)](#), [format\\_with\\_commas\(\)](#), [gene\\_report\(\)](#), [strip\\_plotly\(\)](#)

## Examples

```
d <- example_data()
# QC metrics
qc <- calculate_qc_metrics(d$rnaseq_se)
head(qc)
```

```
# Clinical cleaning
clin <- clean_clinical_data(d$clinical)

# Survival data with cytogenetic markers
surv <- prepare_survival_data(d$clinical, cyto_file = d$cytogenetic)
```

---

`export_h5ad`*Export SummarizedExperiment to H5AD (AnnData) format*

---

### Description

Converts a [SummarizedExperiment::SummarizedExperiment] to AnnData format and writes an H5AD file for Python/scanpy interoperability.

### Usage

```
export_h5ad(se, output_path, assay_name = NULL)
```

### Arguments

<code>se</code>	A [SummarizedExperiment::SummarizedExperiment] object
<code>output_path</code>	Path for the output .h5ad file
<code>assay_name</code>	Which assay to export (default "unstranded", falls back to "counts" then first available)

### Value

The output path (invisibly)

### See Also

Other utilities: [create\\_summary\\_table\(\)](#), [example\\_data\(\)](#), [format\\_file\\_size\(\)](#), [format\\_with\\_commas\(\)](#), [gene\\_report\(\)](#), [strip\\_plotly\(\)](#)

### Examples

```
## Not run:
d <- example_data()
export_h5ad(d$rnaseq_se, "compass.h5ad")

## End(Not run)
```

---

`extract_cytogenetic_data`*Extract cytogenetic markers from clinical data*

---

**Description**

Parses clinical data from GDC to extract translocation and copy number alteration status. GDC clinical data includes FISH results for standard myeloma cytogenetic markers.

**Usage**

```
extract_cytogenetic_data(clinical_dir)
```

**Arguments**

`clinical_dir` Path to clinical data directory (from `download_clinical_data`)

**Value**

Path to saved cytogenetic parquet file

**See Also**

Other cytogenetics: [calculate\\_cooccurrence\(\)](#), [compute\\_riss\(\)](#), [plot\\_cooccurrence\\_heatmap\(\)](#), [plot\\_cytogenetic\\_oncoprint\(\)](#), [plot\\_expression\\_by\\_subtype\(\)](#), [summarize\\_cytogenetics\(\)](#)

**Examples**

```
## Not run:
cyto_file <- extract_cytogenetic_data("data/raw/clinical")
cyto <- arrow::read_parquet(cyto_file)

## End(Not run)
```

---

`extract_risk_table`*Extract risk table from Kaplan-Meier results*

---

**Description**

Generates a number-at-risk table at fixed time points from a `[run_kaplan_meier()]` result. Useful for pre-computing Shiny display data and static API endpoints.

**Usage**

```
extract_risk_table(km_result, times = c(0, 365, 730, 1095, 1460))
```

**Arguments**

km_result	List returned by [run_kaplan_meier()]
times	Numeric vector of time points (in days) at which to evaluate the risk table. Defaults to 'c(0, 365, 730, 1095, 1460)' (0-4 years).

**Value**

Data frame with columns: time, n\_risk, n\_event, n\_censor, survival, and optionally strata (if stratified KM)

**See Also**

Other survival: [plot\\_forest\(\)](#), [plot\\_km\(\)](#), [prepare\\_survival\\_data\(\)](#), [run\\_cox\\_regression\(\)](#), [run\\_kaplan\\_meier\(\)](#), [run\\_km\\_by\\_expression\(\)](#), [run\\_km\\_by\\_markers\(\)](#)

---

filter\_low\_quality      *Filter low-quality samples and genes*

---

**Description**

Filter low-quality samples and genes

**Usage**

```
filter_low_quality(se_data, min_counts = 10, min_samples = 3)
```

**Arguments**

se_data	A SummarizedExperiment object
min_counts	Minimum count threshold per gene
min_samples	Minimum samples with counts above threshold

**Value**

Filtered SummarizedExperiment

---

find\_consensus\_genes *Find consensus DE genes across methods*

---

**Description**

Find consensus DE genes across methods

**Usage**

```
find_consensus_genes(de_results_list, padj_threshold = 0.05, lfc_threshold = 1)
```

**Arguments**

de\_results\_list  
List of DE result objects from different methods

padj\_threshold Adjusted p-value threshold (default: 0.05)

lfc\_threshold Log2 fold change threshold (default: 1)

**Value**

List with consensus gene information

---

format\_file\_size *Format File Size in Human-Readable Format*

---

**Description**

Converts file sizes from bytes to human-readable format with appropriate units

**Usage**

```
format_file_size(size_bytes, digits = 1)
```

**Arguments**

size\_bytes Numeric vector of file sizes in bytes

digits Number of decimal places to show (default: 1)

**Value**

Character vector with formatted file sizes

**See Also**

Other utilities: [create\\_summary\\_table\(\)](#), [example\\_data\(\)](#), [export\\_h5ad\(\)](#), [format\\_with\\_commas\(\)](#), [gene\\_report\(\)](#), [strip\\_plotly\(\)](#)

**Examples**

```
format_file_size(c(1024, 1048576, 5000877192))  
# Returns: "1.0 KB", "1.0 MB", "4.7 GB"
```

---

format_with_commas	<i>Format Number with Thousands Separator</i>
--------------------	---

---

**Description**

Adds commas as thousands separators to large numbers

**Usage**

```
format_with_commas(x)
```

**Arguments**

x                    Numeric value or vector

**Value**

Character vector with formatted numbers

**See Also**

Other utilities: [create\\_summary\\_table\(\)](#), [example\\_data\(\)](#), [export\\_h5ad\(\)](#), [format\\_file\\_size\(\)](#), [gene\\_report\(\)](#), [strip\\_plotly\(\)](#)

**Examples**

```
format_with_commas(1234567)  
# Returns: "1,234,567"
```

---

gene_report	<i>Render a single-gene characterization report</i>
-------------	---

---

**Description**

Renders the parameterized ‘gene-report.qmd’ vignette for a specific gene, producing a self-contained HTML report with expression distribution, survival stratification, cytogenetic context, and gene-gene correlations.

**Usage**

```
gene_report(
  gene,
  output_dir = "results/gene_reports",
  vst_matrix = NULL,
  surv_data = NULL,
  cyto_data = NULL
)
```

**Arguments**

gene	Gene name (e.g., "CD70", "TP53")
output_dir	Directory for the rendered report (default "results/gene_reports/")
vst_matrix	Optional pre-loaded VST matrix (genes x samples). If NULL, the function attempts to load from the targets store.
surv_data	Optional pre-loaded survival data frame. If NULL, attempts to load from the targets store.
cyto_data	Optional pre-loaded cytogenetic data frame

**Value**

Path to the rendered HTML file (invisibly)

**See Also**

Other utilities: [create\\_summary\\_table\(\)](#), [example\\_data\(\)](#), [export\\_h5ad\(\)](#), [format\\_file\\_size\(\)](#), [format\\_with\\_commas\(\)](#), [strip\\_plotly\(\)](#)

**Examples**

```
## Not run:
gene_report("CD70")
gene_report("TP53", output_dir = "results/")

## End(Not run)
```

---

generate\_api\_endpoint *Generate a single API endpoint JSON string*

---

**Description**

Wraps a data frame in standard API envelope with metadata and serializes to JSON. Used by [plan\_api] targets to produce static JSON files.

**Usage**

```
generate_api_endpoint(data, dataset_name, description)
```

**Arguments**

data	A data frame to serialize, or NULL if data is unavailable
dataset_name	Short name for this dataset (e.g. "clinical")
description	Human-readable description

**Value**

A JSON string (character scalar) with 'metadata' and 'data' fields

**See Also**

Other api: [api\\_get\\_clinical\(\)](#), [api\\_get\\_de\\_results\(\)](#), [api\\_get\\_pathways\(\)](#), [api\\_get\\_survival\(\)](#), [api\\_list\\_datasets\(\)](#), [api\\_serve\(\)](#), [generate\\_api\\_index\(\)](#)

**Examples**

```
## Not run:
df <- data.frame(a = 1:3, b = letters[1:3])
json <- generate_api_endpoint(df, "example", "Example data")
cat(json)

## End(Not run)
```

---

generate_api_index	<i>Generate API index metadata</i>
--------------------	------------------------------------

---

**Description**

Creates the index.json content listing all available endpoints, their descriptions, and URLs for the static JSON API.

**Usage**

```
generate_api_index(
  base_url = "https://JohnGavin.github.io/coMmpass-analysis/api/v1"
)
```

**Arguments**

base_url	Base URL for the static API. Defaults to the GitHub Pages URL.
----------	--

**Value**

A list suitable for JSON serialization containing endpoint catalogue, version, and generated timestamp.

**See Also**

Other api: [api\\_get\\_clinical\(\)](#), [api\\_get\\_de\\_results\(\)](#), [api\\_get\\_pathways\(\)](#), [api\\_get\\_survival\(\)](#), [api\\_list\\_datasets\(\)](#), [api\\_serve\(\)](#), [generate\\_api\\_endpoint\(\)](#)

**Examples**

```
## Not run:
idx <- generate_api_index()
cat(jsonlite::toJSON(idx, pretty = TRUE, auto_unbox = TRUE))

## End(Not run)
```

---

```
generate_summary_report
      Generate summary report
```

---

**Description**

Generate summary report

**Usage**

```
generate_summary_report(
  qc_metrics,
  de_genes,
  survival,
  pathways,
  output_dir = "results/reports"
)
```

**Arguments**

qc_metrics	QC metrics data frame
de_genes	DE results with consensus gene information
survival	Survival analysis results
pathways	Pathway analysis results
output_dir	Directory for output report

**Value**

Path to generated report

---

get\_adjustment\_sets    *Get adjustment sets for a given analysis*

---

**Description**

Uses the DAG to identify the minimal sufficient adjustment set for estimating the causal effect of ‘exposure’ on ‘outcome’.

**Usage**

```
get_adjustment_sets(  
  exposure = "cytogenetic_risk",  
  outcome = "overall_survival",  
  dag = NULL  
)
```

**Arguments**

exposure	Character. The exposure variable name in the DAG.
outcome	Character. The outcome variable name in the DAG.
dag	A ‘dagitty’ DAG object. Defaults to [compass_dag()].

**Value**

A list of character vectors, each a minimal sufficient adjustment set.

---

get\_compass\_clinical    *Query GDC for CoMMpass Clinical Data*

---

**Description**

Retrieves clinical data for the MMRF CoMMpass study from the Genomic Data Commons (GDC). This includes patient demographics, disease characteristics, treatment information, and outcomes data.

**Usage**

```
get_compass_clinical()
```

**Value**

A data frame containing clinical data for CoMMpass patients

**See Also**

Other data-acquisition: [acquire\\_commpass\\_data\(\)](#), [download\\_clinical\\_data\(\)](#), [download\\_gdc\\_rnaseq\(\)](#), [download\\_s3\\_subset\(\)](#), [list\\_s3\\_commpass\(\)](#), [query\\_commpass\\_rna\(\)](#)

**Examples**

```
## Not run:
# Get clinical data
clinical <- get_commpass_clinical()

# View first few rows
head(clinical)

# Check available columns
names(clinical)

## End(Not run)
```

---

get\_commpass\_data\_dictionary  
*Get CoMMpass Data Dictionary*

---

**Description**

Returns a tibble documenting all known variables in the CoMMpass dataset, including clinical, biospecimen, and RNA-seq data. Each variable includes its category, data type, units, description, typical range, and a link to the GDC data dictionary.

**Usage**

```
get_commpass_data_dictionary()
```

**Value**

A tibble with columns: variable, category, data\_type, units, description, typical\_range, gdc\_link

**See Also**

Other data-dictionary: [get\\_variable\\_docs\(\)](#)

**Examples**

```
dd <- get_commpass_data_dictionary()
# Filter to clinical variables
dplyr::filter(dd, category == "clinical")
```

---

get_commpass_tbl	<i>Get Lazy DuckDB Table for CoMMpass Data</i>
------------------	--

---

### Description

Returns a lazy 'dplyr::tbl()' backed by DuckDB, reading from parquet files. The connection is managed by the caller and must be disconnected when done.

### Usage

```
get_commpass_tbl(
  data_type = c("clinical", "biospecimen", "rnaseq_counts", "rnaseq_sample_metadata",
    "rnaseq_gene_metadata"),
  data_dir = "data/raw",
  con = NULL
)
```

### Arguments

data_type	One of "clinical", "biospecimen", "rnaseq_counts", "rnaseq_sample_metadata", "rnaseq_gene_metadata"
data_dir	Base directory containing parquet files
con	An existing DuckDB connection. If NULL, a new in-memory connection is created and returned as an attribute of the result.

### Value

A lazy dbplyr tbl. If con was NULL, the DuckDB connection is stored as attr(result, "connection") - caller must disconnect it.

### See Also

Other storage: [query\\_commpass\\_parquet\(\)](#)

### Examples

```
## Not run:
# Create connection and query
con <- DBI::dbConnect(duckdb::duckdb())
clinical_tbl <- get_commpass_tbl("clinical", con = con)

# Chain dplyr operations (lazy - not executed until collect)
result <- clinical_tbl |>
  dplyr::filter(gender == "female") |>
  dplyr::select(submitter_id, age_at_diagnosis, vital_status) |>
  dplyr::collect()

# Clean up
```

```
DBI::dbDisconnect(con, shutdown = TRUE)

## End(Not run)
```

---

get\_counts\_assay      *Get counts assay from SummarizedExperiment*

---

**Description**

GDC STAR-Counts uses 'unstranded', but fallback to 'counts' if present

**Usage**

```
get_counts_assay(se)
```

**Arguments**

se                      SummarizedExperiment object

**Value**

Counts matrix

---

get\_variable\_docs      *Get Extended Documentation for a Variable*

---

**Description**

Returns detailed documentation for a specific variable, including scientific context, calculation methods, and usage notes.

**Usage**

```
get_variable_docs(variable)
```

**Arguments**

variable                Character string naming the variable to document

**Value**

A list with elements: variable, description, scientific\_context, calculation, usage\_notes, references

**See Also**

Other data-dictionary: [get\\_compass\\_data\\_dictionary\(\)](#)

**Examples**

```
docs <- get_variable_docs("age_at_diagnosis")
cat(docs$usage_notes)
```

---

```
integrate_clinical_expression
      Create Integrated Dataset
```

---

**Description**

Combines clinical and expression data with consistent sample IDs

**Usage**

```
integrate_clinical_expression(clinical_clean, expr_clean)
```

**Arguments**

clinical\_clean Cleaned clinical data  
expr\_clean Cleaned expression data

**Value**

List with matched clinical and expression data

**See Also**

Other data-cleaning: [clean\\_clinical\\_data\(\)](#), [clean\\_expression\\_data\(\)](#), [clean\\_treatment\\_data\(\)](#), [summarize\\_treatment\(\)](#)

---

```
list_s3_compass      List AWS S3 CoMMpass Bucket Contents
```

---

**Description**

Lists files available in the public AWS S3 bucket containing MMRF CoMMpass data. The bucket contains RNA-seq, genomic, and clinical data files. This function uses anonymous access to the public bucket.

**Usage**

```
list_s3_compass(prefix = "")
```

**Arguments**

prefix Optional prefix to filter files (e.g., "RNA-seq/", "clinical/")

**Value**

Character vector of S3 object keys (file paths)

**See Also**

Other data-acquisition: [acquire\\_compass\\_data\(\)](#), [download\\_clinical\\_data\(\)](#), [download\\_gdc\\_rnaseq\(\)](#), [download\\_s3\\_subset\(\)](#), [get\\_compass\\_clinical\(\)](#), [query\\_compass\\_rna\(\)](#)

**Examples**

```
## Not run:  
# List all available files (limited to first 100)  
files <- list_s3_compass()  
  
# List only RNA-seq files  
rna_files <- list_s3_compass(prefix = "RNA-seq/")  
  
# Count available files  
length(files)  
  
## End(Not run)
```

---

normalize_rnaseq	<i>Normalize RNA-seq data</i>
------------------	-------------------------------

---

**Description**

Normalize RNA-seq data

**Usage**

```
normalize_rnaseq(se_data, method = "TMM")
```

**Arguments**

se_data	A SummarizedExperiment object with a "counts" assay
method	Normalization method (default: "TMM")

**Value**

SummarizedExperiment with added "logCPM" assay

---

`plot_cooccurrence_heatmap`*Plot co-occurrence heatmap of cytogenetic alterations*

---

**Description**

Displays pairwise co-occurrence patterns as a heatmap. Color indicates  $-\log_{10}(\text{p-value})$  direction: red for co-occurrence, blue for mutual exclusivity.

**Usage**

```
plot_cooccurrence_heatmap(cooccurrence, title = "Cytogenetic Co-occurrence")
```

**Arguments**

<code>cooccurrence</code>	Data frame from <code>calculate_cooccurrence()</code>
<code>title</code>	Plot title

**Value**

A ggplot object

**See Also**

Other cytogenetics: [calculate\\_cooccurrence\(\)](#), [compute\\_riss\(\)](#), [extract\\_cytogenetic\\_data\(\)](#), [plot\\_cytogenetic\\_oncoprint\(\)](#), [plot\\_expression\\_by\\_subtype\(\)](#), [summarize\\_cytogenetics\(\)](#)

**Examples**

```
## Not run:  
cooc <- calculate_cooccurrence(cyto)  
plot_cooccurrence_heatmap(cooc)  
  
## End(Not run)
```

---

`plot_cytogenetic_oncoprint`*Plot cytogenetic oncoprint*

---

**Description**

Creates an oncoprint-style heatmap showing cytogenetic alterations across patients. Rows are alterations, columns are patients. Uses ggplot2 for portability (no ComplexHeatmap dependency).

**Usage**

```
plot_cytogenetic_oncoprint(
  cyto_data,
  markers = NULL,
  sort_by = c("frequency", "risk"),
  title = "Cytogenetic Landscape"
)
```

**Arguments**

cyto_data	Data frame from <code>extract_cytogenetic_data()</code> with columns: <code>patient_id</code> , <code>iss_stage</code> , <code>t_4_14</code> , <code>t_11_14</code> , <code>t_14_16</code> , <code>del_17p</code> , <code>gain_1q</code> , etc.
markers	Character vector of marker columns to include. Default uses all available markers.
sort_by	How to sort patients. One of "frequency" (default) or "risk".
title	Plot title.

**Value**

A ggplot object

**See Also**

Other cytogenetics: [calculate\\_cooccurrence\(\)](#), [compute\\_riss\(\)](#), [extract\\_cytogenetic\\_data\(\)](#), [plot\\_cooccurrence\\_heatmap\(\)](#), [plot\\_expression\\_by\\_subtype\(\)](#), [summarize\\_cytogenetics\(\)](#)

**Examples**

```
## Not run:
cyto <- arrow::read_parquet("data/raw/clinical/cytogenetic_data.parquet")
plot_cytogenetic_oncoprint(cyto)

## End(Not run)
```

---

plot\_dag

---

*Plot the CoMMpass causal DAG*


---

**Description**

Renders the DAG using `ggdag` with sensible defaults.

**Usage**

```
plot_dag(dag = NULL, title = "CoMMpass Causal DAG")
```

**Arguments**

dag            A ‘dagitty’ DAG object. Defaults to [commpass\_dag()].  
title          Plot title.

**Value**

A ggplot object.

---

plot\_enrichment\_barplot  
*Bar plot of enrichment results*

---

**Description**

Creates a horizontal bar plot of the top N enriched pathways, colored by direction (GSEA) or significance (ORA).

**Usage**

```
plot_enrichment_barplot(enrich_result, n = 15L, title = "Enrichment Bar Plot")
```

**Arguments**

enrich\_result    Data frame with pathway enrichment results  
n                Number of top pathways to show (default 15)  
title            Plot title

**Value**

A ggplot2 object

**See Also**

Other pathway: [annotate\\_de\\_results\(\)](#), [annotate\\_genes\(\)](#), [plot\\_enrichment\\_dotplot\(\)](#), [plot\\_gsea\\_running\\_score\(\)](#), [run\\_gsea\(\)](#), [run\\_ora\(\)](#), [run\\_pathway\\_analysis\(\)](#)

plot\_enrichment\_dotplot

*Dot plot of enrichment results*

---

### Description

Creates a dot plot of the top N enriched pathways. Dot size represents gene set size, color represents significance (adjusted p-value or NES).

### Usage

```
plot_enrichment_dotplot(  
  enrich_result,  
  n = 15L,  
  color_by = "padj",  
  title = "Enrichment Dot Plot"  
)
```

### Arguments

enrich_result	Data frame with pathway enrichment results. Must contain columns: pathway, padj, and one of: NES (GSEA), overlap (ORA).
n	Number of top pathways to show (default 15)
color_by	Which metric to color by: "padj" (default) or "NES"
title	Plot title

### Value

A ggplot2 object

### See Also

Other pathway: [annotate\\_de\\_results\(\)](#), [annotate\\_genes\(\)](#), [plot\\_enrichment\\_barplot\(\)](#), [plot\\_gsea\\_running\\_score\(\)](#), [run\\_gsea\(\)](#), [run\\_ora\(\)](#), [run\\_pathway\\_analysis\(\)](#)

---

plot\_expression\_by\_subtype

*Plot gene expression by cytogenetic subtype*

---

### Description

Creates violin + box plots showing gene expression stratified by cytogenetic marker status (positive vs negative). Adds Wilcoxon or t-test p-value annotation.

**Usage**

```
plot_expression_by_subtype(
  expr_matrix,
  cyto_data,
  gene,
  markers = NULL,
  test = c("wilcox", "t.test")
)
```

**Arguments**

expr_matrix	Numeric matrix (genes x samples) of transformed expression values (e.g., VST). Column names should match 'cyto_data\$patient_id'.
cyto_data	Data frame with 'patient_id' and marker columns (values: "positive"/"negative"/NA)
gene	Gene name (rowname of 'expr_matrix')
markers	Character vector of marker columns. Default auto-detects.
test	Statistical test: "wilcox" (default) or "t.test"

**Value**

A ggplot object with faceted violin/box plots, one panel per marker

**See Also**

Other cytogenetics: [calculate\\_cooccurrence\(\)](#), [compute\\_riss\(\)](#), [extract\\_cytogenetic\\_data\(\)](#), [plot\\_cooccurrence\\_heatmap\(\)](#), [plot\\_cytogenetic\\_oncprint\(\)](#), [summarize\\_cytogenetics\(\)](#)

---

plot_forest	<i>Plot forest plot of Cox hazard ratios</i>
-------------	--

---

**Description**

Creates a forest plot from Cox regression results, showing hazard ratios with 95

**Usage**

```
plot_forest(cox_result, title = "Cox Regression Forest Plot")
```

**Arguments**

cox_result	List from [run_cox_regression()] with 'hazard_ratios', 'concordance', 'n', 'n_events'
title	Plot title

**Value**

A ggplot object

**See Also**

Other survival: [extract\\_risk\\_table\(\)](#), [plot\\_km\(\)](#), [prepare\\_survival\\_data\(\)](#), [run\\_cox\\_regression\(\)](#), [run\\_kaplan\\_meier\(\)](#), [run\\_km\\_by\\_expression\(\)](#), [run\\_km\\_by\\_markers\(\)](#)

---

plot\_gene\_correlation *Scatter plot of gene-gene correlation*

---

**Description**

Creates a scatter plot with regression line, correlation coefficient, and p-value annotation from a `[correlate_genes()]` result.

**Usage**

```
plot_gene_correlation(cor_result, title = NULL)
```

**Arguments**

<code>cor_result</code>	List from <code>[correlate_genes()]</code>
<code>title</code>	Plot title (default auto-generated)

**Value**

A ggplot object

**See Also**

Other gene-correlation: [correlate\\_genes\(\)](#), [correlate\\_genes\\_batch\(\)](#)

---

plot\_gsea\_running\_score  
*GSEA running enrichment score plot*

---

**Description**

Shows the running enrichment score for a specific gene set from fgsea results, along with the ranked gene positions.

**Usage**

```
plot_gsea_running_score(  
  gsea_result,  
  gene_set_name,  
  gene_sets = NULL,  
  title = NULL  
)
```

**Arguments**

gsea_result	List from run_gsea() containing results and ranked_genes
gene_set_name	Name of the gene set to plot
gene_sets	Named list of gene sets (required for tick marks). If NULL, attempts to retrieve from msigdb using the collection.
title	Plot title (default: gene set name)

**Value**

A ggplot2 object

**See Also**

Other pathway: [annotate\\_de\\_results\(\)](#), [annotate\\_genes\(\)](#), [plot\\_enrichment\\_barplot\(\)](#), [plot\\_enrichment\\_dotplot\(\)](#), [run\\_gsea\(\)](#), [run\\_ora\(\)](#), [run\\_pathway\\_analysis\(\)](#)

---

plot_heatmap_de	<i>Heatmap of top DE genes</i>
-----------------	--------------------------------

---

**Description**

Creates a heatmap of the top differentially expressed genes using z-score scaled expression values.

**Usage**

```
plot_heatmap_de(
  expr_matrix,
  de_results,
  n_genes = 50L,
  annotation_df = NULL,
  title = "Top DE Genes"
)
```

**Arguments**

expr_matrix	Numeric matrix of transformed expression (genes x samples)
de_results	Data frame with DE results (rownames = gene IDs)
n_genes	Number of top DE genes to show (default 50)
annotation_df	Optional data frame of sample annotations for column annotation
title	Plot title

**Value**

A ggplot2 object (tile-based heatmap)

**See Also**

Other differential-expression: [compute\\_pca\(\)](#), [plot\\_ma\(\)](#), [plot\\_pca\(\)](#), [plot\\_volcano\(\)](#), [run\\_deseq2\(\)](#), [run\\_vst\(\)](#), [summarize\\_de\\_methods\(\)](#)

---

plot\_km

*Plot Kaplan-Meier curve*

---

**Description**

Creates a ggplot2-based KM survival curve from a survfit object. Supports stratified and unstratified curves with optional log-rank p-value.

**Usage**

```
plot_km(  
  km_result,  
  title = "Kaplan-Meier Survival Curve",  
  xlab = "Time (days)",  
  time_breaks = NULL  
)
```

**Arguments**

km_result	List from [run_kaplan_meier()] with 'fit', 'logrank_p', 'median_survival', 'n_per_group', 'strata'
title	Plot title
xlab	X-axis label
time_breaks	Sequence of time breaks for x-axis (in days). Default marks every 365 days.

**Value**

A ggplot object

**See Also**

Other survival: [extract\\_risk\\_table\(\)](#), [plot\\_forest\(\)](#), [prepare\\_survival\\_data\(\)](#), [run\\_cox\\_regression\(\)](#), [run\\_kaplan\\_meier\(\)](#), [run\\_km\\_by\\_expression\(\)](#), [run\\_km\\_by\\_markers\(\)](#)

---

plot_ma	<i>MA plot of DE results</i>
---------	------------------------------

---

**Description**

Plots mean expression (x-axis) vs log2 fold change (y-axis), a standard diagnostic for RNA-seq DE analysis.

**Usage**

```
plot_ma(  
  de_results,  
  lfc_threshold = 1,  
  padj_threshold = 0.05,  
  title = "MA Plot"  
)
```

**Arguments**

de_results	Data frame with DE results
lfc_threshold	Log2 fold-change threshold for coloring (default 1)
padj_threshold	Adjusted p-value threshold (default 0.05)
title	Plot title

**Value**

A ggplot2 object

**See Also**

Other differential-expression: [compute\\_pca\(\)](#), [plot\\_heatmap\\_de\(\)](#), [plot\\_pca\(\)](#), [plot\\_volcano\(\)](#), [run\\_deseq2\(\)](#), [run\\_vst\(\)](#), [summarize\\_de\\_methods\(\)](#)

---

plot_pca	<i>PCA plot</i>
----------	-----------------

---

**Description**

Creates a PCA scatter plot from pre-computed PCA data.

**Usage**

```
plot_pca(pca_data, color_by = NULL, shape_by = NULL, title = "PCA")
```

**Arguments**

pca_data	List from compute_pca() with coords and var_explained
color_by	Column name in coords to color by (default NULL)
shape_by	Column name in coords to use for shape (default NULL)
title	Plot title

**Value**

A ggplot2 object

**See Also**

Other differential-expression: [compute\\_pca\(\)](#), [plot\\_heatmap\\_de\(\)](#), [plot\\_ma\(\)](#), [plot\\_volcano\(\)](#), [run\\_deseq2\(\)](#), [run\\_vst\(\)](#), [summarize\\_de\\_methods\(\)](#)

---

plot_volcano	<i>Volcano plot of DE results</i>
--------------	-----------------------------------

---

**Description**

Creates a ggplot2 volcano plot with significance thresholds and optional gene labels.

**Usage**

```
plot_volcano(
  de_results,
  lfc_threshold = 1,
  padj_threshold = 0.05,
  n_label = 10L,
  title = "Volcano Plot"
)
```

**Arguments**

de_results	Data frame with DE results (must have padj and log2FC columns – auto-detected from DESeq2/edgeR/limma conventions)
lfc_threshold	Log2 fold-change threshold for significance (default 1)
padj_threshold	Adjusted p-value threshold (default 0.05)
n_label	Number of top genes to label (default 10)
title	Plot title

**Value**

A ggplot2 object

**See Also**

Other differential-expression: [compute\\_pca\(\)](#), [plot\\_heatmap\\_de\(\)](#), [plot\\_ma\(\)](#), [plot\\_pca\(\)](#), [run\\_deseq2\(\)](#), [run\\_vst\(\)](#), [summarize\\_de\\_methods\(\)](#)

---

prepare\_survival\_data *Prepare survival data from clinical and cytogenetic data*

---

**Description**

Constructs a data frame suitable for survival analysis from GDC clinical data. Uses ‘days\_to\_death’ for deceased patients and ‘days\_to\_last\_follow\_up’ for censored (alive) patients. Merges with cytogenetic risk groups when available.

**Usage**

```
prepare_survival_data(clinical_data, cyto_file = NULL)
```

**Arguments**

`clinical_data` Cleaned clinical data frame (from `clean_clinical_data`)

`cyto_file` Path to cytogenetic parquet file (from `extract_cytogenetic_data`), or `NULL` to skip cytogenetic integration

**Value**

Data frame with columns: `patient_id`, `time_days`, `status` (0=censored, 1=dead), `age_years`, `gender`, `iss_stage`, `risk_group`, plus individual cytogenetic markers

**See Also**

Other survival: [extract\\_risk\\_table\(\)](#), [plot\\_forest\(\)](#), [plot\\_km\(\)](#), [run\\_cox\\_regression\(\)](#), [run\\_kaplan\\_meier\(\)](#), [run\\_km\\_by\\_expression\(\)](#), [run\\_km\\_by\\_markers\(\)](#)

---

query\_commpass\_parquet

*Query CoMMpass Parquet Files*

---

**Description**

Opens an in-memory DuckDB connection, creates a VIEW on the specified parquet file, optionally applies filters, and returns a data frame.

**Usage**

```
query_commpass_parquet(
  data_type = c("clinical", "biospecimen", "rnaseq_counts", "rnaseq_sample_metadata",
    "rnaseq_gene_metadata"),
  data_dir = "data/raw",
  filters = NULL,
  collect = TRUE
)
```

**Arguments**

data_type	One of "clinical", "biospecimen", "rnaseq_counts", "rnaseq_sample_metadata", "rnaseq_gene_metadata"
data_dir	Base directory containing parquet files
filters	Optional named list of filters. Names are column names, values are vectors of allowed values (used in WHERE ... IN (...))
collect	If TRUE (default), collect results into a data frame. If FALSE, return a lazy tbl for further dplyr operations.

**Value**

A data frame (if collect=TRUE) or lazy dbplyr tbl

**See Also**

Other storage: [get\\_commpass\\_tbl\(\)](#)

**Examples**

```
## Not run:
# Read all clinical data
clinical <- query_commpass_parquet("clinical")

# Filter to specific patients
subset <- query_commpass_parquet(
  "clinical",
  filters = list(gender = "female", vital_status = "Alive")
)

# Get lazy tbl for chaining
tbl <- query_commpass_parquet("clinical", collect = FALSE)
result <- tbl |> dplyr::filter(gender == "female") |> dplyr::collect()

## End(Not run)
```

---

query_commpass_rna	<i>Query GDC for CoMMpass RNA-seq Metadata</i>
--------------------	--

---

**Description**

Queries the Genomic Data Commons (GDC) API for Multiple Myeloma Research Foundation (MMRF) CoMMpass study RNA-seq data. This function returns a query object that can be used with other TCGAbiolinks functions to download and prepare the data.

**Usage**

```
query_commpass_rna()
```

**Value**

A GDCquery object containing metadata for RNA-seq samples

**See Also**

Other data-acquisition: [acquire\\_commpass\\_data\(\)](#), [download\\_clinical\\_data\(\)](#), [download\\_gdc\\_rnaseq\(\)](#), [download\\_s3\\_subset\(\)](#), [get\\_commpass\\_clinical\(\)](#), [list\\_s3\\_commpass\(\)](#)

**Examples**

```
## Not run:  
# Query for RNA-seq data  
query <- query_commpass_rna()  
  
# Use the query to download data  
# GDCdownload(query)  
# data <- GDCprepare(query)  
  
## End(Not run)
```

---

render_de_report	<i>Render DE analysis report</i>
------------------	----------------------------------

---

**Description**

Render DE analysis report

**Usage**

```
render_de_report(de_results, output_dir = "results/reports")
```

**Arguments**

de\_results      DE results object  
output\_dir      Directory for output report

**Value**

Path to generated report

---

run\_cox\_regression      *Run Cox proportional hazards regression*

---

**Description**

Fits a Cox PH model with specified covariates. Returns the model object, hazard ratios with confidence intervals, and concordance index.

**Usage**

```
run_cox_regression(surv_data, covariates = c("age_years", "gender"))
```

**Arguments**

surv\_data      Data frame from [prepare\_survival\_data()]  
covariates      Character vector of covariate column names

**Value**

List with components: - 'model': coxph object - 'hazard\_ratios': data frame with HR, CI, p-values - 'concordance': C-index - 'n': number of patients in model - 'n\_events': number of events - 'ph\_test': cox.zph result for PH assumption

**See Also**

Other survival: [extract\\_risk\\_table\(\)](#), [plot\\_forest\(\)](#), [plot\\_km\(\)](#), [prepare\\_survival\\_data\(\)](#), [run\\_kaplan\\_meier\(\)](#), [run\\_km\\_by\\_expression\(\)](#), [run\\_km\\_by\\_markers\(\)](#)

---

`run_deseq2`*Run DESeq2 differential expression analysis*

---

## Description

Runs DESeq2 with optional apeglm log-fold-change shrinkage and optional paired longitudinal design for within-patient comparisons.

## Usage

```
run_deseq2(  
  se_data,  
  clinical_data,  
  design_formula = ~condition,  
  shrink_lfc = TRUE,  
  paired = FALSE  
)
```

## Arguments

<code>se_data</code>	A SummarizedExperiment object
<code>clinical_data</code>	Clinical data frame
<code>design_formula</code>	Design formula for the model
<code>shrink_lfc</code>	If TRUE, apply apeglm LFC shrinkage (default TRUE)
<code>paired</code>	If TRUE, use paired longitudinal design with patient as blocking factor. Requires 'patient_id' and 'visit' columns in clinical_data. Only patients with $\geq 2$ timepoints are included.

## Value

List with DE results including raw and shrunken (if requested)

## See Also

Other differential-expression: [compute\\_pca\(\)](#), [plot\\_heatmap\\_de\(\)](#), [plot\\_ma\(\)](#), [plot\\_pca\(\)](#), [plot\\_volcano\(\)](#), [run\\_vst\(\)](#), [summarize\\_de\\_methods\(\)](#)

---

run_edger	<i>Run edgeR differential expression analysis</i>
-----------	---

---

**Description**

Run edgeR differential expression analysis

**Usage**

```
run_edger(se_data, clinical_data, design_formula = ~condition)
```

**Arguments**

se\_data            A SummarizedExperiment object  
clinical\_data    Clinical data frame  
design\_formula    Design formula for the model

**Value**

List with DE results

---

run_gsea	<i>Run Gene Set Enrichment Analysis (GSEA)</i>
----------	--

---

**Description**

Performs pre-ranked GSEA using fgsea on differential expression results. Genes are ranked by their test statistic (DESeq2 Wald stat or  $-\log_{10}(\text{pvalue}) * \text{sign}(\log_2\text{FC})$ ). Gene sets come from MSigDB via msigdb.

**Usage**

```
run_gsea(  
  de_results,  
  gene_sets = "hallmark",  
  gene_id_type = "ensembl_gene",  
  min_size = 15L,  
  max_size = 500L  
)
```

**Arguments**

<code>de_results</code>	Data frame of DE results. Must contain gene identifiers (as rownames or in a <code>gene/gene_id</code> column) and at least one of: <code>stat</code> (DESeq2 Wald statistic), or <code>log2FoldChange/logFC + pvalue/PValue</code> .
<code>gene_sets</code>	Character string specifying the MSigDB collection: "hallmark" (default), "kegg", "reactome", "go_bp", "go_mf", "go_cc", "c2", "c7". Alternatively, a named list of character vectors (custom gene sets).
<code>gene_id_type</code>	Type of gene identifiers: "ensembl_gene" (default), "gene_symbol", or "entrez_gene".
<code>min_size</code>	Minimum gene set size (default: 15).
<code>max_size</code>	Maximum gene set size (default: 500).

**Value**

List with components:

**results** Data frame with pathway, NES, pval, padj, size, leadingEdge columns.

**n\_gene\_sets** Number of gene sets tested.

**n\_significant** Number significant at  $\text{padj} < 0.05$ .

**top\_gene\_sets** Top 20 enriched gene sets as a data frame.

**ranked\_genes** Named numeric vector of gene ranks used.

**collection** Gene set collection used.

**See Also**

Other pathway: [annotate\\_de\\_results\(\)](#), [annotate\\_genes\(\)](#), [plot\\_enrichment\\_barplot\(\)](#), [plot\\_enrichment\\_dotplot\(\)](#), [plot\\_gsea\\_running\\_score\(\)](#), [run\\_ora\(\)](#), [run\\_pathway\\_analysis\(\)](#)

---

<code>run_kaplan_meier</code>	<i>Run Kaplan-Meier analysis</i>
-------------------------------	----------------------------------

---

**Description**

Fits a Kaplan-Meier survival curve, optionally stratified by a grouping variable. Returns the survfit object, log-rank test, and summary statistics.

**Usage**

```
run_kaplan_meier(surv_data, strata = NULL)
```

**Arguments**

<code>surv_data</code>	Data frame from [ <a href="#">prepare_survival_data()</a> ] with columns 'time_days' and 'status'
<code>strata</code>	Column name to stratify by (e.g. "risk_group", "iss_stage", "gender"), or NULL for overall curve

**Value**

List with components: - 'fit': survfit object - 'logrank': survdiff result (NULL if no strata) - 'logrank\_p': log-rank p-value (NA if no strata) - 'median\_survival': named vector of median survival times - 'n\_per\_group': sample sizes per stratum - 'strata': the stratification variable used

**See Also**

Other survival: [extract\\_risk\\_table\(\)](#), [plot\\_forest\(\)](#), [plot\\_km\(\)](#), [prepare\\_survival\\_data\(\)](#), [run\\_cox\\_regression\(\)](#), [run\\_km\\_by\\_expression\(\)](#), [run\\_km\\_by\\_markers\(\)](#)

---

run\_km\_by\_expression *Run KM analysis stratified by gene expression level*

---

**Description**

Splits patients into groups based on expression of a single gene (e.g., median split) and runs Kaplan-Meier analysis with log-rank test. Connects differential expression results to clinical outcomes.

**Usage**

```
run_km_by_expression(
  surv_data,
  expr_matrix,
  gene,
  split = c("median", "tertile", "quartile", "top_bottom_20"),
  min_per_group = 5L
)
```

**Arguments**

surv_data	Data frame from [ <a href="#">prepare_survival_data()</a> ] with columns 'patient_id', 'time_days', 'status'
expr_matrix	Numeric matrix (genes x samples) of transformed expression values (e.g., VST). Column names must match 'surv_data\$patient_id'.
gene	Gene name (must be a rowname of 'expr_matrix')
split	Split method: "median" (default), "tertile", "quartile", or "top_bottom_20"
min_per_group	Minimum patients per group (default 5)

**Value**

List compatible with [[plot\\_km\(\)](#)]: fit, logrank\_p, median\_survival, n\_per\_group, strata, gene, split\_method. Returns a list with 'fit = NULL' and a 'note' if requirements not met.

**See Also**

Other survival: [extract\\_risk\\_table\(\)](#), [plot\\_forest\(\)](#), [plot\\_km\(\)](#), [prepare\\_survival\\_data\(\)](#), [run\\_cox\\_regression\(\)](#), [run\\_kaplan\\_meier\(\)](#), [run\\_km\\_by\\_markers\(\)](#)

---

run_km_by_markers	<i>Run KM analysis for each individual cytogenetic marker</i>
-------------------	---

---

**Description**

Iterates over cytogenetic markers and runs stratified KM analysis for each. Markers with fewer than 'min\_positive' positive cases are skipped.

**Usage**

```
run_km_by_markers(surv_data, markers = NULL, min_positive = 3L)
```

**Arguments**

surv_data	Data frame from [prepare_survival_data()]
markers	Character vector of marker column names. Default auto-detects.
min_positive	Minimum number of positive cases to run analysis (default 3)

**Value**

Named list of KM results, one per marker

**See Also**

Other survival: [extract\\_risk\\_table\(\)](#), [plot\\_forest\(\)](#), [plot\\_km\(\)](#), [prepare\\_survival\\_data\(\)](#), [run\\_cox\\_regression\(\)](#), [run\\_kaplan\\_meier\(\)](#), [run\\_km\\_by\\_expression\(\)](#)

---

run_limma	<i>Run limma differential expression analysis</i>
-----------	---

---

**Description**

Run limma differential expression analysis

**Usage**

```
run_limma(se_data, clinical_data, design_formula = ~condition)
```

**Arguments**

se_data	A SummarizedExperiment object
clinical_data	Clinical data frame
design_formula	Design formula for the model

**Value**

List with DE results

---

run\_ora *Run Over-Representation Analysis (ORA)*

---

### Description

Tests whether a set of significant genes is over-represented in gene set collections using Fisher's exact test (hypergeometric distribution). Gene sets come from MSigDB via msigdb.

### Usage

```
run_ora(
  sig_genes,
  universe,
  gene_sets = "hallmark",
  gene_id_type = "ensembl_gene",
  min_size = 10L,
  max_size = 500L
)
```

### Arguments

sig_genes	Character vector of significant gene identifiers.
universe	Character vector of all tested gene identifiers (background).
gene_sets	Character string specifying the MSigDB collection (same options as [run_gsea()]), or a named list of character vectors.
gene_id_type	Type of gene identifiers: "ensembl_gene" (default), "gene_symbol", or "entrez_gene".
min_size	Minimum gene set size (default: 10).
max_size	Maximum gene set size (default: 500).

### Value

List with components:

**results** Data frame with pathway, overlap, gene\_set\_size, universe\_size, p\_value, padj, odds\_ratio, overlapping\_genes.

**n\_pathways\_tested** Number of pathways tested.

**n\_pathways\_enriched** Number significant at padj < 0.05.

**n\_sig\_genes** Number of input significant genes.

**top\_pathways** Top 20 enriched pathways as a data frame.

### See Also

Other pathway: [annotate\\_de\\_results\(\)](#), [annotate\\_genes\(\)](#), [plot\\_enrichment\\_barplot\(\)](#), [plot\\_enrichment\\_dotplot\(\)](#), [plot\\_gsea\\_running\\_score\(\)](#), [run\\_gsea\(\)](#), [run\\_pathway\\_analysis\(\)](#)

---

run\_pathway\_analysis *Run pathway enrichment analysis*

---

**Description**

Wrapper that runs ORA on consensus DE genes. Uses MSigDB Hallmark gene sets by default.

**Usage**

```
run_pathway_analysis(de_genes, method = "hallmark")
```

**Arguments**

de_genes	Consensus DE result from [find_consensus_genes()] with consensus_genes character vector and by_method list.
method	Gene set collection: "hallmark" (default), "kegg", "reactome", "go_bp".

**Value**

List with ORA results (see [run\_ora()]).

**See Also**

Other pathway: [annotate\\_de\\_results\(\)](#), [annotate\\_genes\(\)](#), [plot\\_enrichment\\_barplot\(\)](#), [plot\\_enrichment\\_dotplot\(\)](#), [plot\\_gsea\\_running\\_score\(\)](#), [run\\_gsea\(\)](#), [run\\_ora\(\)](#)

---

run\_vst *Run variance stabilizing transformation*

---

**Description**

Wrapper around DESeq2::vst() for use in visualization. VST is appropriate for PCA, heatmaps, and clustering – not for DE testing (which uses raw counts).

**Usage**

```
run_vst(se_data, blind = TRUE)
```

**Arguments**

se_data	SummarizedExperiment with "counts" assay
blind	Logical, whether VST should be blind to the design (default TRUE for exploratory analysis)

**Value**

SummarizedExperiment with "vst" assay added, or NULL if DESeq2 is not available

**See Also**

Other differential-expression: [compute\\_pca\(\)](#), [plot\\_heatmap\\_de\(\)](#), [plot\\_ma\(\)](#), [plot\\_pca\(\)](#), [plot\\_volcano\(\)](#), [run\\_deseq2\(\)](#), [summarize\\_de\\_methods\(\)](#)

---

save_timestamped	<i>Save results with timestamp</i>
------------------	------------------------------------

---

**Description**

Save results with timestamp

**Usage**

```
save_timestamped(object, base_name, dir = "results")
```

**Arguments**

object	R object to save
base_name	Base name for the output file
dir	Directory to save to (default: "results")

---

setup_logging	<i>Setup logging</i>
---------------	----------------------

---

**Description**

Setup logging

**Usage**

```
setup_logging(log_file = NULL)
```

**Arguments**

log_file	Optional path to a log file. If NULL, logs to console.
----------	--

---

strip_plotly	<i>Strip plotly closure bloat for compact serialization</i>
--------------	---

---

**Description**

Plotly htmlwidgets capture parent environments in closures, causing 60-80x size inflation when serialized via saveRDS/targets. This function runs plotly\_build() to resolve lazy data, then replaces closure environments with emptyenv() to eliminate the bloat.

**Usage**

```
strip_plotly(p)
```

**Arguments**

p                    A plotly htmlwidget object

**Value**

The same plotly object with closures stripped (much smaller on disk)

**See Also**

Other utilities: [create\\_summary\\_table\(\)](#), [example\\_data\(\)](#), [export\\_h5ad\(\)](#), [format\\_file\\_size\(\)](#), [format\\_with\\_commas\(\)](#), [gene\\_report\(\)](#)

**Examples**

```
## Not run:
p <- plotly::plot_ly(mtcars, x = ~mpg, type = "histogram")
object.size(p)        # large
p2 <- strip_plotly(p)
object.size(p2)      # small

## End(Not run)
```

---

summarize_cytogenetics	<i>Summarize cytogenetic alteration frequencies</i>
------------------------	---

---

**Description**

Computes frequency and percentage for each marker and risk group.

**Usage**

```
summarize_cytogenetics(cyto_data)
```

**Arguments**

cyto\_data      Data frame from extract\_cytogenetic\_data()

**Value**

Data frame with marker, n\_positive, n\_tested, pct columns

**See Also**

Other cytogenetics: [calculate\\_cooccurrence\(\)](#), [compute\\_riss\(\)](#), [extract\\_cytogenetic\\_data\(\)](#), [plot\\_cooccurrence\\_heatmap\(\)](#), [plot\\_cytogenetic\\_oncoprint\(\)](#), [plot\\_expression\\_by\\_subtype\(\)](#)

**Examples**

```
## Not run:  
cyto <- arrow::read_parquet("data/raw/clinical/cytogenetic_data.parquet")  
summarize_cytogenetics(cyto)  
  
## End(Not run)
```

---

summarize\_data      *Generate summary statistics*

---

**Description**

Generate summary statistics

**Usage**

```
summarize_data(se_data)
```

**Arguments**

se\_data      A SummarizedExperiment object with a "counts" assay

---

summarize\_de\_methods *Summarize DE results across methods*

---

**Description**

Creates a summary table showing the number of significant genes per DE method.

**Usage**

```
summarize_de_methods(de_list, padj_threshold = 0.05, lfc_threshold = 1)
```

**Arguments**

`de_list` Named list of DE result lists (each with `results_table`)  
`padj_threshold` Adjusted p-value threshold (default 0.05)  
`lfc_threshold` Log2 fold-change threshold (default 1)

**Value**

Data frame with `method`, `n_tested`, `n_sig`, `n_up`, `n_down`

**See Also**

Other differential-expression: [compute\\_pca\(\)](#), [plot\\_heatmap\\_de\(\)](#), [plot\\_ma\(\)](#), [plot\\_pca\(\)](#), [plot\\_volcano\(\)](#), [run\\_deseq2\(\)](#), [run\\_vst\(\)](#)

---

summarize\_treatment *Summarize treatment lines per patient*

---

**Description**

Creates a one-row-per-patient summary of treatment history.

**Usage**

```
summarize_treatment(treatment_clean)
```

**Arguments**

`treatment_clean`  
Cleaned treatment data from `[clean_treatment_data()]`

**Value**

Data frame with columns: `patient_id`, `n_lines`, `first_regimen`, `first_regimen_class`, `had_transplant`, `best_overall_response`

**See Also**

Other data-cleaning: [clean\\_clinical\\_data\(\)](#), [clean\\_expression\\_data\(\)](#), [clean\\_treatment\\_data\(\)](#), [integrate\\_clinical\\_expression\(\)](#)

# Index

- \* **api**
  - api\_get\_clinical, 6
  - api\_get\_de\_results, 7
  - api\_get\_pathways, 8
  - api\_get\_survival, 8
  - api\_list\_datasets, 9
  - api\_serve, 9
  - generate\_api\_endpoint, 28
  - generate\_api\_index, 29
- \* **cytogenetics**
  - calculate\_cooccurrence, 10
  - compute\_riss, 15
  - extract\_cytogenetic\_data, 24
  - plot\_cooccurrence\_heatmap, 37
  - plot\_cytogenetic\_oncoprint, 37
  - plot\_expression\_by\_subtype, 40
  - summarize\_cytogenetics, 59
- \* **data-acquisition**
  - acquire\_compass\_data, 4
  - download\_clinical\_data, 19
  - download\_gdc\_rnaseq, 20
  - download\_s3\_subset, 21
  - get\_compass\_clinical, 31
  - list\_s3\_compass, 35
  - query\_compass\_rna, 49
- \* **data-cleaning**
  - clean\_clinical\_data, 12
  - clean\_expression\_data, 13
  - clean\_treatment\_data, 13
  - integrate\_clinical\_expression, 35
  - summarize\_treatment, 61
- \* **data-dictionary**
  - get\_compass\_data\_dictionary, 32
  - get\_variable\_docs, 34
- \* **differential-expression**
  - compute\_pca, 14
  - correlate\_genes, 16
  - correlate\_genes\_batch, 16
  - plot\_gene\_correlation, 42
  - plot\_heatmap\_de, 43
  - plot\_ma, 45
  - plot\_pca, 45
  - plot\_volcano, 46
  - run\_deseq2, 51
  - run\_vst, 57
  - summarize\_de\_methods, 61
- \* **gene-correlation**
  - correlate\_genes, 16
  - correlate\_genes\_batch, 16
  - plot\_gene\_correlation, 42
- \* **pathway**
  - annotate\_de\_results, 5
  - annotate\_genes, 6
  - plot\_enrichment\_barplot, 39
  - plot\_enrichment\_dotplot, 40
  - plot\_gsea\_running\_score, 42
  - run\_gsea, 52
  - run\_ora, 56
  - run\_pathway\_analysis, 57
- \* **quality-control**
  - calculate\_qc\_metrics, 11
- \* **storage**
  - get\_compass\_tbl, 33
  - query\_compass\_parquet, 47
- \* **survival**
  - extract\_risk\_table, 24
  - plot\_forest, 41
  - plot\_km, 44
  - prepare\_survival\_data, 47
  - run\_cox\_regression, 50
  - run\_kaplan\_meier, 53
  - run\_km\_by\_expression, 54
  - run\_km\_by\_markers, 55
- \* **utilities**
  - create\_summary\_table, 17
  - example\_data, 22
  - export\_h5ad, 23
  - format\_file\_size, 26

- format\_with\_commas, 27
  - gene\_report, 27
  - get\_counts\_assay, 34
  - strip\_plotly, 59
- acquire\_compass\_data, 4, 19–21, 32, 36, 49
- annotate\_de\_results, 5, 6, 39, 40, 43, 53, 56, 57
- annotate\_genes, 5, 6, 39, 40, 43, 53, 56, 57
- api\_get\_clinical, 6, 7–10, 29, 30
- api\_get\_de\_results, 7, 7–10, 29, 30
- api\_get\_pathways, 7, 8, 9, 10, 29, 30
- api\_get\_survival, 7, 8, 8–10, 29, 30
- api\_list\_datasets, 7, 8, 9, 9, 10, 29, 30
- api\_serve, 7, 8, 9, 9, 29, 30
- calculate\_cooccurrence, 10, 15, 24, 37, 38, 41, 60
- calculate\_qc\_metrics, 11
- check\_adjustment, 11
- check\_dependencies, 12
- clean\_clinical\_data, 12, 13, 14, 35, 62
- clean\_expression\_data, 12, 13, 14, 35, 62
- clean\_treatment\_data, 12, 13, 13, 35, 62
- compass\_dag, 14
- compute\_pca, 14, 44–47, 51, 58, 61
- compute\_riss, 11, 15, 24, 37, 38, 41, 60
- correlate\_genes, 16, 17, 42
- correlate\_genes\_batch, 16, 16, 42
- create\_project\_dirs, 17
- create\_summary\_table, 17, 22, 23, 26–28, 59
- download\_aws\_data, 18
- download\_clinical\_data, 4, 19, 20, 21, 32, 36, 49
- download\_gdc\_rnaseq, 4, 19, 20, 21, 32, 36, 49
- download\_s3\_subset, 4, 19, 20, 21, 32, 36, 49
- example\_data, 18, 22, 23, 26–28, 59
- export\_h5ad, 18, 22, 23, 26–28, 59
- extract\_cytogenetic\_data, 11, 15, 24, 37, 38, 41, 60
- extract\_risk\_table, 24, 42, 44, 47, 50, 54, 55
- filter\_low\_quality, 25
- find\_consensus\_genes, 26
- format\_file\_size, 18, 22, 23, 26, 27, 28, 59
- format\_with\_commas, 18, 22, 23, 26, 27, 28, 59
- gene\_report, 18, 22, 23, 26, 27, 27, 59
- generate\_api\_endpoint, 7–10, 28, 30
- generate\_api\_index, 7–10, 29, 29
- generate\_summary\_report, 30
- get\_adjustment\_sets, 31
- get\_compass\_clinical, 4, 19–21, 31, 36, 49
- get\_compass\_data\_dictionary, 32, 34
- get\_compass\_tbl, 33, 48
- get\_counts\_assay, 34
- get\_variable\_docs, 32, 34
- integrate\_clinical\_expression, 12–14, 35, 62
- list\_s3\_compass, 4, 19–21, 32, 35, 49
- normalize\_rnaseq, 36
- plot\_cooccurrence\_heatmap, 11, 15, 24, 37, 38, 41, 60
- plot\_cytogenetic\_oncoprint, 11, 15, 24, 37, 37, 41, 60
- plot\_dag, 38
- plot\_enrichment\_barplot, 5, 6, 39, 40, 43, 53, 56, 57
- plot\_enrichment\_dotplot, 5, 6, 39, 40, 43, 53, 56, 57
- plot\_expression\_by\_subtype, 11, 15, 24, 37, 38, 40, 60
- plot\_forest, 25, 41, 44, 47, 50, 54, 55
- plot\_gene\_correlation, 16, 17, 42
- plot\_gsea\_running\_score, 5, 6, 39, 40, 42, 53, 56, 57
- plot\_heatmap\_de, 14, 43, 45–47, 51, 58, 61
- plot\_km, 25, 42, 44, 47, 50, 54, 55
- plot\_ma, 14, 44, 45, 46, 47, 51, 58, 61
- plot\_pca, 14, 44, 45, 45, 47, 51, 58, 61
- plot\_volcano, 14, 44, 45, 46, 46, 51, 58, 61
- prepare\_survival\_data, 25, 42, 44, 47, 50, 54, 55
- query\_compass\_parquet, 33, 47
- query\_compass\_rna, 4, 19–21, 32, 36, 49

render\_de\_report, 49  
run\_cox\_regression, 25, 42, 44, 47, 50, 54,  
55  
run\_deseq2, 14, 44–47, 51, 58, 61  
run\_edger, 52  
run\_gsea, 5, 6, 39, 40, 43, 52, 56, 57  
run\_kaplan\_meier, 25, 42, 44, 47, 50, 53, 54,  
55  
run\_km\_by\_expression, 25, 42, 44, 47, 50,  
54, 54, 55  
run\_km\_by\_markers, 25, 42, 44, 47, 50, 54, 55  
run\_limma, 55  
run\_ora, 5, 6, 39, 40, 43, 53, 56, 57  
run\_pathway\_analysis, 5, 6, 39, 40, 43, 53,  
56, 57  
run\_vst, 14, 44–47, 51, 57, 61  
  
save\_timestamped, 58  
setup\_logging, 58  
strip\_plotly, 18, 22, 23, 26–28, 59  
summarize\_cytogenetics, 11, 15, 24, 37, 38,  
41, 59  
summarize\_data, 60  
summarize\_de\_methods, 14, 44–47, 51, 58,  
61  
summarize\_treatment, 12–14, 35, 61