

# COSMIC 67

*Julian Gehring, EMBL Heidelberg*

October 30, 2021

## Contents

1	Introduction . . . . .	1
2	Accessing and Using the Data . . . . .	1
3	Data Provenance . . . . .	4
3.1	COSMIC Mutations . . . . .	4
3.2	Cancer Gene Census. . . . .	4
4	Data Source. . . . .	5
5	References . . . . .	5
6	Session Info. . . . .	5

## 1 Introduction

---

The *COSMIC.67* package provides the curated mutations published with the COSMIC release version 67 (2013-10-24). Both variants found in coding and non-coding regions are included and offered as a single object of class 'CollapsedVCF' and a bgzipped and tabix-index 'VCF' file.

Additionally, the package contains the Cancer Gene Census, a list of genes causally linked to cancer.

## 2 Accessing and Using the Data

---

`library(VariantAnnotation)`

*Loading required package: BiocGenerics*

*Attaching package: 'BiocGenerics'*

*The following objects are masked from 'package:stats':*

*IQR, mad, sd, var, xtabs*

*The following objects are masked from 'package:base':*

## COSMIC 67

*Filter, Find, Map, Position, Reduce, anyDuplicated, append, as.data.frame, basename, cbind, colnames, dirname, do.call, duplicated, eval, evalq, get, grep, grepl, intersect, is.unsorted, lapply, mapply, match, mget, order, paste, pmax, pmax.int, pmin, pmin.int, rank, rbind, rownames, sapply, setdiff, sort, table, tapply, union, unique, unsplit, which.max, which.min*

Loading required package: *MatrixGenerics*

Loading required package: *matrixStats*

Attaching package: *'MatrixGenerics'*

The following objects are masked from *'package:matrixStats'*:

*colAlls, colAnyNAs, colAnys, colAvgPerRowSet, colCollapse, colCounts, colCummaxs, colCummins, colCumprods, colCumsums, colDiffs, colIQRDiffs, colIQRs, colLogSumExps, colMadDiffs, colMads, colMaxs, colMeans2, colMedians, colMins, colOrderStats, colProds, colQuantiles, colRanges, colRanks, colSdDiffs, colSds, colSums2, colTabulates, colVarDiffs, colVars, colWeightedMads, colWeightedMeans, colWeightedMedians, colWeightedSds, colWeightedVars, rowAlls, rowAnyNAs, rowAnys, rowAvgPerColSet, rowCollapse, rowCounts, rowCummaxs, rowCummins, rowCumprods, rowCumsums, rowDiffs, rowIQRDiffs, rowIQRs, rowLogSumExps, rowMadDiffs, rowMads, rowMaxs, rowMeans2, rowMedians, rowMins, rowOrderStats, rowProds, rowQuantiles, rowRanges, rowRanks, rowSdDiffs, rowSds, rowSums2, rowTabulates, rowVarDiffs, rowVars, rowWeightedMads, rowWeightedMeans, rowWeightedMedians, rowWeightedSds, rowWeightedVars*

Loading required package: *GenomeInfoDb*

Loading required package: *S4Vectors*

Loading required package: *stats4*

Attaching package: *'S4Vectors'*

The following objects are masked from *'package:base'*:

*I, expand.grid, unname*

Loading required package: *IRanges*

Loading required package: *GenomicRanges*

Loading required package: *SummarizedExperiment*

Loading required package: *Biobase*

Welcome to *Bioconductor*

## COSMIC 67

Vignettes contain introductory material; view with `'browseVignettes()'`. To cite Bioconductor, see `'citation("Biobase")'`, and for packages `'citation("pkgname")'`.

Attaching package: `'Biobase'`

The following object is masked from `'package:MatrixGenerics'`:

`rowMedians`

The following objects are masked from `'package:matrixStats'`:

`anyMissing`, `rowMedians`

Loading required package: `Rsamtools`

Loading required package: `Biostrings`

Loading required package: `XVector`

Attaching package: `'Biostrings'`

The following object is masked from `'package:base'`:

`strsplit`

Attaching package: `'VariantAnnotation'`

The following object is masked from `'package:base'`:

`tabulate`

`library(GenomicRanges)`

`data(package = "COSMIC.67")`

`data(cosmic_67, package = "COSMIC.67")`

`tp53_range = GRanges("17", IRanges(7565097, 7590856))`

`vcf_path = system.file("vcf", "cosmic_67.vcf.gz", package = "COSMIC.67")`

`cosmic_tp53 = readVcf(vcf_path, genome = "GRCh37", ScanVcfParam(which = tp53_range))`

`cosmic_tp53`

`class: CollapsedVCF`

`dim: 5892 0`

`rowRanges(vcf):`

`GRanges with 5 metadata columns: paramRangeID, REF, ALT, QUAL, FILTER`

`info(vcf):`

`DataFrame with 5 columns: GENE, STRAND, CDS, AA, CNT`

`info(header(vcf)):`

	Number	Type	Description
GENE	1	String	Gene name
STRAND	1	String	Gene strand
CDS	1	String	CDS annotation
AA	1	String	Peptide annotation

```
CNT      1      Integer How many samples have this mutation
geno(vcf):
List of length 0:
```

```
data(cgc_67, package = "COSMIC.67")
head(cgc_67)
```

	SYMBOL	ENTREZID	ENSEMBL
1	ABI1	10006	ENSG000000136754
2	ABL1	25	ENSG000000097007
3	ABL2	27	ENSG000000143322
4	ACSL3	2181	ENSG000000123983
5	CASC5	57082	ENSG000000137812
6	MLLT11	10962	ENSG000000213190

For details on the collection and curation of the original data, please see the webpage of the COSMIC project: <http://cancer.sanger.ac.uk/cancergenome/projects/cosmic/>.

## 3 Data Provenance

---

### 3.1 COSMIC Mutations

The following steps are performed for importing and processing of the VCF data:

1. Downloading of the VCF files 'CosmicCodingMuts\_v67\_20131024.vcf.gz' and 'Cosmic-NonCodingVariants\_v67\_20131024.vcf.gz' from 'ftp://ngs.sanger.ac.uk/production/cosmic/' to 'inst/raw/
2. Importing of both files to R using 'readVcf'.
3. Sorting of the seqlevels and adding 'seqinfo' data for the toplevel chromosomes of 'GRCh37'.
4. Merging of both objects, sorting according to genomic position.
5. Converting the object to class `VariantAnnotation::VRanges`.
6. Converting the 'character' columns to 'factors'.
7. Saving the merged object to 'data/cosmic\_v67\_vcf.rda'.
8. Exporting the merged object as a bgzipped and tabix-indexed 'VCF' to 'inst/vcf/cosmic\_v67.vcf.gz'.

### 3.2 Cancer Gene Census

The following steps are performed for importing and processing of the Cancer Gene Census data:

1. Downloading of the 'cancer\_gene\_census.tsv' file from [ftp://ftp.sanger.ac.uk/pub/CGP/cosmic/data\\_export](ftp://ftp.sanger.ac.uk/pub/CGP/cosmic/data_export) to 'inst/raw'.
2. Import of the files as a data frame.
3. Annotation of the 'HGNC' and 'ENSEMBLID' identifiers, using the 'ENTREZ gene ID' as query with the 'org.Hs.eg.db' object.

4. Saving the object to 'data/cgc\_67.rda'.

## 4 Data Source

---

The mutation data was obtained from the Sanger Institute Catalogue Of Somatic Mutations In Cancer web site, <http://www.sanger.ac.uk/cosmic>

Bamford et al (2004):

The COSMIC (Catalogue of Somatic Mutations in Cancer) database and website.

Br J Cancer, 91,355-358.

For details on the usage and redistribution of the data, please see [ftp://ftp.sanger.ac.uk/pub/CGP/cosmic/GUIDELINES\\_ON\\_THE\\_USE\\_OF\\_THIS\\_DATA.txt](ftp://ftp.sanger.ac.uk/pub/CGP/cosmic/GUIDELINES_ON_THE_USE_OF_THIS_DATA.txt).

## 5 References

---

- <http://cancer.sanger.ac.uk/cancergenome/projects/cosmic/>
- [http://nar.oxfordjournals.org/content/39/suppl\\_1/D945.long](http://nar.oxfordjournals.org/content/39/suppl_1/D945.long)
- [ftp://ftp.sanger.ac.uk/pub/CGP/cosmic/GUIDELINES\\_ON\\_THE\\_USE\\_OF\\_THIS\\_DATA.txt](ftp://ftp.sanger.ac.uk/pub/CGP/cosmic/GUIDELINES_ON_THE_USE_OF_THIS_DATA.txt)

## 6 Session Info

---

R version 4.1.1 (2021-08-10)

Platform: x86\_64-pc-linux-gnu (64-bit)

Running under: Ubuntu 20.04.3 LTS

Matrix products: default

BLAS: /home/biocbuild/bbs-3.14-bioc/R/lib/libRblas.so

LAPACK: /home/biocbuild/bbs-3.14-bioc/R/lib/libRlapack.so

locale:

```
[1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
[3] LC_TIME=en_GB             LC_COLLATE=C
[5] LC_MONETARY=en_US.UTF-8   LC_MESSAGES=en_US.UTF-8
[7] LC_PAPER=en_US.UTF-8     LC_NAME=C
[9] LC_ADDRESS=C              LC_TELEPHONE=C
[11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
```

attached base packages:

```
[1] stats4      stats      graphics  grDevices  utils      datasets
[7] methods     base
```

other attached packages:

```
[1] VariantAnnotation_1.40.0  Rsamtools_2.10.0
[3] Biostrings_2.62.0        XVector_0.34.0
[5] SummarizedExperiment_1.24.0 Biobase_2.54.0
[7] GenomicRanges_1.46.0     GenomeInfoDb_1.30.0
```

## COSMIC 67

```
[9] IRanges_2.28.0          S4Vectors_0.32.0
[11] MatrixGenerics_1.6.0    matrixStats_0.61.0
[13] BiocGenerics_0.40.0     knitr_1.36
```

loaded via a namespace (and not attached):

```
[1] httr_1.4.2              bit64_4.0.5
[3] assertthat_0.2.1       highr_0.9
[5] BiocManager_1.30.16    BiocFileCache_2.2.0
[7] blob_1.2.2             BSgenome_1.62.0
[9] GenomeInfoDbData_1.2.7 yaml_2.2.1
[11] progress_1.2.2         pillar_1.6.4
[13] RSQLite_2.2.8          lattice_0.20-45
[15] glue_1.4.2             digest_0.6.28
[17] htmltools_0.5.2       Matrix_1.3-4
[19] XML_3.99-0.8           pkgconfig_2.0.3
[21] biomaRt_2.50.0         zlibbioc_1.40.0
[23] purrr_0.3.4           BiocParallel_1.28.0
[25] tibble_3.1.5          KEGGREST_1.34.0
[27] generics_0.1.1        ellipsis_0.3.2
[29] cachem_1.0.6          GenomicFeatures_1.46.1
[31] magrittr_2.0.1        crayon_1.4.1
[33] memoise_2.0.0         evaluate_0.14
[35] fansi_0.5.0           xml2_1.3.2
[37] tools_4.1.1           prettyunits_1.1.1
[39] hms_1.1.1             BiocStyle_2.22.0
[41] BiocIO_1.4.0          lifecycle_1.0.1
[43] stringr_1.4.0         DelayedArray_0.20.0
[45] AnnotationDbi_1.56.1  compiler_4.1.1
[47] rlang_0.4.12          grid_4.1.1
[49] RCurl_1.98-1.5        rjson_0.2.20
[51] rappdirs_0.3.3       bitops_1.0-7
[53] rmarkdown_2.11       restfulr_0.0.13
[55] DBI_1.1.1             curl_4.3.2
[57] R6_2.5.1             GenomicAlignments_1.30.0
[59] dplyr_1.0.7          rtracklayer_1.54.0
[61] fastmap_1.1.0         bit_4.0.4
[63] utf8_1.2.2           filelock_1.0.2
[65] stringi_1.7.5        parallel_4.1.1
[67] Rcpp_1.0.7           vctrs_0.3.8
[69] png_0.1-7            dbplyr_2.1.1
[71] tidyselect_1.1.1     xfun_0.27
```