

# COSMIC 67

*Julian Gehring, EMBL Heidelberg*

October 29, 2020

## 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 . . . . .           | 5 |
| 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
```

```
Loading required package: parallel
```

```
Attaching package: 'BiocGenerics'
```

## COSMIC 67

The following objects are masked from 'package:parallel':

*clusterApply, clusterApplyLB, clusterCall,  
clusterEvalQ, clusterExport, clusterMap, parApply,  
parCapply, parLapply, parLapplyLB, parRapply,  
parSapply, parSapplyLB*

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

*IQR, mad, sd, var, xtabs*

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

*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'

## COSMIC 67

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

`expand.grid`

Loading required package: `IRanges`

Loading required package: `GenomicRanges`

Loading required package: `SummarizedExperiment`

Loading required package: `Biobase`

Welcome to Bioconductor

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
```

## COSMIC 67

```
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    | ENSG00000136754 |
| 2 | ABL1   | 25       | ENSG00000097007 |
| 3 | ABL2   | 27       | ENSG00000143322 |
| 4 | ACSL3  | 2181     | ENSG00000123983 |
| 5 | CASC5  | 57082    | ENSG00000137812 |
| 6 | MLLT11 | 10962    | ENSG00000213190 |

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 4.0.3 (2020-10-10)

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

Running under: Ubuntu 18.04.5 LTS

Matrix products: default

BLAS: /home/biocbuild/bbs-3.12-bioc/R/lib/libRblas.so

LAPACK: /home/biocbuild/bbs-3.12-bioc/R/lib/libRlapack.so

locale:

```
[1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
[3] LC_TIME=en_US.UTF-8      LC_COLLATE=C
```

## COSMIC 67

```
[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 parallel stats graphics grDevices utils
[7] datasets methods base
```

other attached packages:

```
[1] VariantAnnotation_1.36.0 Rsamtools_2.6.0
[3] Biostrings_2.58.0 XVector_0.30.0
[5] SummarizedExperiment_1.20.0 Biobase_2.50.0
[7] GenomicRanges_1.42.0 GenomeInfoDb_1.26.0
[9] IRanges_2.24.0 S4Vectors_0.28.0
[11] MatrixGenerics_1.2.0 matrixStats_0.57.0
[13] BiocGenerics_0.36.0 knitr_1.30
```

loaded via a namespace (and not attached):

```
[1] Rcpp_1.0.5 lattice_0.20-41
[3] prettyunits_1.1.1 assertthat_0.2.1
[5] digest_0.6.27 BiocFileCache_1.14.0
[7] R6_2.5.0 RSQLite_2.2.1
[9] evaluate_0.14 highr_0.8
[11] httr_1.4.2 pillar_1.4.6
[13] zlibbioc_1.36.0 rlang_0.4.8
[15] GenomicFeatures_1.42.0 progress_1.2.2
[17] curl_4.3 blob_1.2.1
[19] Matrix_1.2-18 rmarkdown_2.5
[21] BiocParallel_1.24.0 stringr_1.4.0
[23] RCurl_1.98-1.2 bit_4.0.4
[25] biomaRt_2.46.0 DelayedArray_0.16.0
[27] rtracklayer_1.50.0 compiler_4.0.3
[29] xfun_0.18 pkgconfig_2.0.3
[31] askpass_1.1 htmltools_0.5.0
[33] tidyselect_1.1.0 openssl_1.4.3
[35] tibble_3.0.4 GenomeInfoDbData_1.2.4
[37] XML_3.99-0.5 crayon_1.3.4
[39] dplyr_1.0.2 dbplyr_1.4.4
[41] GenomicAlignments_1.26.0 rappdirs_0.3.1
[43] bitops_1.0-6 grid_4.0.3
[45] lifecycle_0.2.0 DBI_1.1.0
[47] magrittr_1.5 stringi_1.5.3
[49] xml2_1.3.2 ellipsis_0.3.1
[51] vctrs_0.3.4 generics_0.0.2
[53] BiocStyle_2.18.0 tools_4.0.3
[55] bit64_4.0.5 BSgenome_1.58.0
[57] glue_1.4.2 purrr_0.3.4
[59] hms_0.5.3 yaml_2.2.1
[61] AnnotationDbi_1.52.0 BiocManager_1.30.10
[63] memoise_1.1.0
```