

# Package ‘RTCGA’

April 23, 2016

**Title** The Cancer Genome Atlas Data Integration

**Version** 1.0.2

**Date** 2015-01-18

**Author** Marcin Kosinski <m.p.kosinski@gmail.com>, Przemyslaw Biecek  
<przemyslaw.biecek@gmail.com>

**Maintainer** Marcin Kosinski <m.p.kosinski@gmail.com>

**Description** The Cancer Genome Atlas (TCGA) Data Portal provides a platform for researchers to search, download, and analyze data sets generated by TCGA. It contains clinical information, genomic characterization data, and high level sequence analysis of the tumor genomes. The key is to understand genomics to improve cancer care. RTCGA package offers download and integration of the variety and volume of TCGA data using patient barcode key, what enables easier data possession. This may have an beneficial influence on impact on development of science and improvement of patients' treatment. Furthermore, RTCGA package transforms TCGA data to tidy form which is convenient to use.

**BugReports** <https://github.com/RTCGA/RTCGA/issues>

**License** GPL-2

**LazyLoad** yes

**LazyData** yes

**Depends** R (>= 3.2.0), knitr

**Imports** XML, assertthat, stringi, rvest, data.table, magrittr, xml2

**Suggests** testthat, pander, RTCGA.rnaseq, RTCGA.clinical,  
RTCGA.mutations

**Repository** Bioconductor

**biocViews** Software, DataImport, DataRepresentation, Preprocessing

**VignetteBuilder** knitr

**NeedsCompilation** no

**RoxygenNote** 5.0.1

## R topics documented:

RTCGA-package . . . . .	2
checkTCGA . . . . .	3
datasetsTCGA . . . . .	4
downloadTCGA . . . . .	6
infoTCGA . . . . .	7
readTCGA . . . . .	8

<b>Index</b>	<b>15</b>
--------------	-----------

---

RTCGA-package	<i>The Cancer Genome Atlas data integration</i>
---------------	---

---

### Description

The Cancer Genome Atlas (TCGA) Data Portal provides a platform for researchers to search, download, and analyze data sets generated by TCGA. It contains clinical information, genomic characterization data, and high level sequence analysis of the tumor genomes. The key is to understand genomics to improve cancer care. RTCGA package offers download and integration of the variety and volume of TCGA data using patient barcode key, what enables easier data possession. This may have an beneficial influence on impact on development of science and improvement of patients' treatment. Furthermore, RTCGA package transforms TCGA data to form which is convenient to use in R statistical package. Those data transformations can be a part of statistical analysis pipeline which can be more reproducible with RTCGA

### Details

For more detailed information visit **RTCGA** wiki on [Github](#).

### Author(s)

Marcin Kosinski [aut, cre] < m.p.kosinski@gmail.com >  
 Przemyslaw Biecek [aut] < przemyslaw.biecek@gmail.com >

### See Also

Other RTCGA: [checkTCGA](#), [datasetsTCGA](#), [downloadTCGA](#), [infoTCGA](#), [readTCGA](#)

### Examples

```
## Not run:
browseVignettes('RTCGA')

## End(Not run)
```

## Description

The checkTCGA function let's to check

- DataSets: TCGA datasets' names for current release date and cohort.
- Dates: TCGA datasets' dates of release.

## Usage

```
checkTCGA(what, cancerType, date = NULL)
```

## Arguments

what	One of DataSets or Dates.
cancerType	A character of length 1 containing abbreviation (Cohort code - <a href="http://gdac.broadinstitute.org/">http://gdac.broadinstitute.org/</a> ) of types of cancers to check for.
date	A NULL or character specifying from which date informations should be checked. By default (date = NULL) the newest available date is used. All available dates can be checked on <a href="http://gdac.broadinstitute.org/runs/">http://gdac.broadinstitute.org/runs/</a> or by using checkTCGA('Dates') function. Required format 'YYYY-MM-DD'.

## Details

- If what='DataSets' enables to check TCGA datasets' names for current release date and cohort.
- If what='Dates' enables to check dates of TCGA datasets' releases.

## Value

- If what='DataSets' a data.frame of available datasets' names (to pass to the [downloadTCGA](#) function) and sizes.
- If what='Dates' a vector of available dates to pass to the [downloadTCGA](#) function.

## See Also

Other RTCGA: [RTCGA-package](#), [datasetsTCGA](#), [downloadTCGA](#), [infoTCGA](#), [readTCGA](#)

## Examples

```
#####

# names for current release date and cohort
checkTCGA('DataSets', 'BRCA' )
## Not run:
checkTCGA('DataSets', 'OV', tail(checkTCGA('Dates'))[3] )
#checkTCGA('DataSets', 'OV', checkTCGA('Dates')[5] ) # error

## End(Not run)
# dates of TCGA datasets' releases.
checkTCGA('Dates')

#####
## Not run:
# TCGA datasets' names availability for
# current release date and cancer type.

releaseDate <- '2015-08-21'
cancerTypes <- c('OV', 'BRCA')

cancerTypes %>% sapply(function(element){
  grep(x = checkTCGA('DataSets', element, releaseDate)[, 1],
    pattern = 'humanmethylation450', value = TRUE) %>%
    as.vector()
})

## End(Not run)
```

---

datasetsTCGA

*RTCGA.data - The Family of R Packages with Data from The Cancer Genome Atlas Study*

---

## Description

Snapshots of the clinical, mutations, CNVs, rnaseq, RPPA, mRNA, miRNASeq and methylation datasets from the 2015-11-01 release date (check all dates of release with `checkTCGA('Dates')`) are included in the `RTCGA.data` family (factory) that contains 9 packages:

- **RTCGA.rnaseq** [rnaseq](#)
- **RTCGA.clinical** [clinical](#)
- **RTCGA.mutations** [mutations](#)
- **RTCGA.CNV**
- **RTCGA.RPPA**
- **RTCGA.mRNA**

- **RTCGA.miRNASeq**
- **RTCGA.methylation**
- **RTCGA.PANCAN12** (not from TCGA)

### Details

For more detailed information visit **RTCGA.data** website <https://rtcga.github.io/RTCGA>. One can install all data packages with `installTCGA (devel)`.

### Author(s)

Marcin Kosinski [aut, cre] < m.p.kosinski@gmail.com >  
Przemyslaw Biecek [aut] < przemyslaw.biecek@gmail.com >  
Witold Chodor [aut] < witoldchodor@gmail.com >

### See Also

Other RTCGA: [RTCGA-package](#), [checkTCGA](#), [downloadTCGA](#), [infoTCGA](#), [readTCGA](#)

### Examples

```
# installation of packages containing snapshots
# of TCGA project's datasets

## Not run:

## RTCGA GitHub development newest versions
library(RTCGA)
?installTCGA

## Bioconductor releases
source('http://bioconductor.org/biocLite.R')
biocLite(RTCGA.clinical)
biocLite(RTCGA.mutations)
biocLite(RTCGA.rnaseq)
biocLite(RTCGA.CNV)
biocLite(RTCGA.RPPA)
biocLite(RTCGA.mRNA)
biocLite(RTCGA.miRNASeq)
biocLite(RTCGA.methylation)

# use cases and examples + more data info
browseVignettes('RTCGA')

## End(Not run)
```

---

 downloadTCGA

*Download TCGA data*


---

### Description

Enables to download TCGA data from specified dates of releases of concrete Cohorts of cancer types. Pass a name of required dataset to the `dataSet` parameter. By default the Merged Clinical `dataSet` is downloaded (value `dataSet = 'Merge_Clinical.Level_1'`) from the newest available date of the release.

### Usage

```
downloadTCGA(cancerTypes, dataSet = "Merge_Clinical.Level_1", destDir,
             date = NULL, untarFile = TRUE, removeTar = TRUE, allDataSets = FALSE)
```

### Arguments

<code>cancerTypes</code>	A character vector containing abbreviations (Cohort code) of types of cancers to download from <a href="http://gdac.broadinstitute.org/">http://gdac.broadinstitute.org/</a> . For easy access from R check details below.
<code>dataSet</code>	A part of the name of <code>dataSet</code> to be downloaded from <a href="http://gdac.broadinstitute.org/runs/">http://gdac.broadinstitute.org/runs/</a> . By default the Merged Clinical <code>dataSet</code> is downloaded (value <code>dataSet = 'Merge_Clinical.Level_1'</code> ). Available datasets' names can be checked using <a href="#">checkTCGA</a> function.
<code>destDir</code>	A character specifying a directory into which <code>dataSets</code> will be downloaded.
<code>date</code>	A NULL or character specifying from which date <code>dataSets</code> should be downloaded. By default ( <code>date = NULL</code> ) the newest available date is used. All available dates can be checked on <a href="http://gdac.broadinstitute.org/runs/">http://gdac.broadinstitute.org/runs/</a> or by using <a href="#">checkTCGA</a> function. Required format 'YYYY-MM-DD'.
<code>untarFile</code>	Logical - should the downloaded file be untarred. Default is TRUE.
<code>removeTar</code>	Logical - should the downloaded .tar file be removed after untarring. Default is TRUE.
<code>allDataSets</code>	Logical - should download all datasets matching <code>dataSet</code> parameter or only the first one (without FFPE phrase if possible).

### Details

All cohort names can be checked using: `sub( x = names( infoTCGA() ), '-counts', '' )`.

### Value

No values. It only downloads files.

### See Also

Other R/TCGA: [RTCGA-package](#), [checkTCGA](#), [datasetsTCGA](#), [infoTCGA](#), [readTCGA](#)

## Examples

```
dir.create( 'hre')

downloadTCGA( cancerTypes = 'ACC', dataSet = 'miR_gene_expression',
  destDir = 'hre', date = tail( checkTCGA('Dates'), 2 )[1] )

## Not run:
downloadTCGA( cancerTypes = c('BRCA', 'OV'), destDir = 'hre',
  date = tail( checkTCGA('Dates'), 2 )[1] )

## End(Not run)
```

---

infoTCGA

*Information about cohorts from TCGA project*

---

## Description

Function restores codes and counts for each cohort from TCGA project.

## Usage

```
infoTCGA()
```

## Value

A list with a tabular information from <http://gdac.broadinstitute.org/>.

## See Also

Other RTCGA: [RTCGA-package](#), [checkTCGA](#), [datasetsTCGA](#), [downloadTCGA](#), [readTCGA](#)

## Examples

```
infoTCGA()

(cohorts <- infoTCGA()) %>%
  rownames() %>%
  sub('-counts', '', x=.)
```

---

readTCGA	<i>Read TCGA data to the tidy format</i>
----------	--

---

## Description

readTCGA function allows to read unzipped files:

- clinical data - Merge\_Clinical.Level\_1
- rnaseq data (genes' expressions) - rnaseqv2\_\_illuminahiseq\_rnaseqv2
- genes' mutations data - Mutation\_Packager\_Calls.Level
- Reverse phase protein array data (RPPA) - protein\_normalization\_\_data.Level\_3
- Merge transcriptome agilent data (mRNA) - Merge\_transcriptome\_\_agilentg4502a\_07\_3\_\_unc\_edu\_\_Level\_3\_\_RSEM
- miRNASeq data - Merge\_mirnaseq\_\_illuminaga\_mirnaseq\_\_bcgsc\_ca\_\_Level\_3\_\_miR\_gene\_expression\_\_data.Level\_3 or "Merge\_mirnaseq\_\_illuminahiseq\_mirnaseq\_\_bcgsc\_ca\_\_Level\_3\_\_miR\_gene\_expression\_\_data.Level\_3"
- methylation data - Merge\_methylation\_\_humanmethylation27
- isoforms data - Merge\_rnaseqv2\_\_illuminahiseq\_rnaseqv2\_\_unc\_edu\_\_Level\_3\_\_RSEM\_isoforms\_normalized

from TCGA project. Those files can be easily downloaded with [downloadTCGA](#) function. See examples.

## Usage

```
readTCGA(path, dataType, ...)
```

## Arguments

path	See details and examples.
dataType	One of 'clinical', 'rnaseq', 'mutations', 'RPPA', 'mRNA', 'miRNASeq', 'methylation', 'isoforms', depending on which type of data user is trying to read in the tidy format.
...	Further arguments passed to the <a href="#">as.data.frame</a> .

## Details

All cohort names can be checked using: `sub( x = names( infoTCGA() ), '-counts', '' )`.

Parameter path specification:

- If dataType = 'clinical' a path to a cancerType.clin.merged.txt file.
- If dataType = 'mutations' a path to the unzipped folder Mutation\_Packager\_Calls.Level containing .maf files.
- If dataType = 'rnaseq' a path to the unzipped file rnaseqv2\_\_illuminahiseq\_rnaseqv2\_\_unc\_edu\_\_Level\_3\_\_RSEM
- If dataType = 'RPPA' a path to the unzipped file in folder protein\_normalization\_\_data.Level\_3.
- If dataType = 'mRNA' a path to the unzipped file cancerType.transcriptome\_\_agilentg4502a\_07\_3\_\_unc\_edu\_\_L
- If dataType = 'miRNASeq' a path to unzipped files cancerType.mirnaseq\_\_illuminahiseq\_mirnaseq\_\_bcgsc\_ca\_\_ or cancerType.mirnaseq\_\_illuminaga\_mirnaseq\_\_bcgsc\_ca\_\_Level\_3\_\_miR\_gene\_expression\_\_data.data.t
- If dataType = 'methylation' a path to unzipped files cancerType.methylation\_\_humanmethylation27\_\_jhuusc
- If dataType = 'isoforms' a path to unzipped files cancerType.rnaseqv2\_\_illuminahiseq\_rnaseqv2\_\_unc\_edu\_\_



**Value**

An output:

- If `dataType = 'clinical'` a `data.frame` with clinical data.
- If `dataType = 'rnaseq'` a `data.frame` with rnaseq data.
- If `dataType = 'mutations'` a `data.frame` with mutations data.
- If `dataType = 'RPPA'` a `data.frame` with RPPA data.
- If `dataType = 'mRNA'` a `data.frame` with mRNA data.
- If `dataType = 'miRNASeq'` a `data.frame` with miRNASeq data.
- If `dataType = 'methylation'` a `data.frame` with methylation data.
- If `dataType = 'isoforms'` a `data.frame` with isoforms data.

**Author(s)**

Marcin Kosinski, <m.p.kosinski@gmail.com>

Witold Chodor, <witoldchodor@gmail.com>

**See Also**

Other R/TCGA: [RTCGA-package](#), [checkTCGA](#), [datasetsTCGA](#), [downloadTCGA](#), [infoTCGA](#)

**Examples**

```
## Not run:

#####
##### clinical
#####

dir.create('data')

# downloading clinical data
# dataset = "clinical" is default parameter so we may omit it
downloadTCGA( cancerTypes = c('BRCA', 'OV'),
              destDir = 'data' )

# reading datasets
sapply( c('BRCA', 'OV'), function( element ){
  folder <- grep( paste0( '(_',element,'\\.', '|','_',element,'-FFPE)', '.*Clinical' ),
                 list.files('data/'),value = TRUE)
  path <- paste0( 'data/', folder, '/', element, '.clin.merged.txt' )
  assign( value = readTCGA( path, 'clinical' ),
          x = paste0(element, '.clin.data'), envir = .GlobalEnv)
})

#####
##### rnaseq
```

```
#####

dir.create('data2')

# downloading rnaseq data
downloadTCGA( cancerTypes = 'BRCA',
              dataSet = 'rnaseqv2__illuminahisec_rnaseqv2__unc_edu__Level_3__RSEM_genes_normalized__data.Level1',
              destDir = 'data2' )

# shortening paths and directories
list.files( 'data2/' ) %>%
  file.path( 'data2', . ) %>%
  file.rename( to = substr(.,start=1,stop=50))

# reading data
list.files( 'data2/' ) %>%
  file.path( 'data2', . ) -> folder

folder %>%
  list.files %>%
  file.path( folder, . ) %>%
  grep( pattern = 'illuminahisec', x = ., value = TRUE) -> pathRNA
readTCGA( path = pathRNA, dataType = 'rnaseq' ) -> my_data

#####
##### mutations
#####

# Example directory in which untarred data will be stored
dir.create('data3')

downloadTCGA( cancerTypes = 'OV',
              dataSet = 'Mutation_Packager_Calls.Level1',
              destDir = 'data3' )

# reading data
list.files( 'data3/' ) %>%
  file.path( 'data3', . ) -> folder

readTCGA(folder, 'mutations') -> mut_file

#####
##### methylation
#####

# Example directory in which untarred data will be stored
dir.create('data4')

# Download KIRP methylation data and store it in data4 folder
cancerType = "KIRP"
downloadTCGA(cancerTypes = cancerType,
```

```

        dataSet = "Merge_methylation__humanmethylation27",
        destDir = "data4")

# Shorten path of subdirectory with KIRP methylation data
list.files(path = "data4", full.names = TRUE) %>%
  file.rename(from = ., to = file.path("data4", paste0(cancerType, ".methylation")))

# Remove manifest.txt file
list.files(path = "data4", full.names = TRUE) %>%
  list.files(path = ., full.names = TRUE) %>%
  grep("MANIFEST.txt", x = ., value = TRUE) %>%
  file.remove()

# Read KIRP methylation data
path <- list.files(path = "data4", full.names = TRUE) %>%
  list.files(path = ., full.names = TRUE)

KIRP.methylation <- readTCGA(path, dataType = "methylation")

#####
##### RPPA
#####

# Directory in which untarred data will be stored
dir.create('data5')

# Download BRCA RPPA data and store it in data5 folder
cancerType = "BRCA"
downloadTCGA(cancerTypes = cancerType,
             dataSet = "protein_normalization__data.Level_3",
             destDir = "data5")

# Shorten path of subdirectory with BRCA RPPA data
list.files(path = "data5", full.names = TRUE) %>%
  file.rename(from = ., to = file.path("data5", paste0(cancerType, ".RPPA")))

# Remove manifest.txt file
list.files(path = "data5", full.names = TRUE) %>%
  list.files(path = ., full.names = TRUE) %>%
  grep("MANIFEST.txt", x = ., value = TRUE) %>%
  file.remove()

# Read BRCA RPPA data
path <- list.files(path = "data5", full.names = TRUE) %>%
  list.files(path = ., full.names = TRUE)

BRCA.RPPA <- readTCGA(path, dataType = "RPPA")

#####
##### mRNA
#####

```

```

# Directory in which untarred data will be stored
dir.create('data6')

# Download UCEC mRNA data and store it in data6 folder
cancerType = "UCEC"
downloadTCGA(cancerTypes = cancerType,
             dataSet = "Merge_transcriptome__agilentg4502a_07_3__unc_edu__Level_3__unc_lowess_normalization_gene_
             destDir = "data6")

# Shorten path of subdirectory with UCEC mRNA data
list.files(path = "data6", full.names = TRUE) %>%
  file.rename(from = ., to = file.path("data6",paste0(cancerType, ".mRNA")))

# Remove manifest.txt file
list.files(path = "data6", full.names = TRUE) %>%
  list.files(path = ., full.names = TRUE) %>%
  grep("MANIFEST.txt", x = ., value = TRUE) %>%
  file.remove()

# Read UCEC mRNA data
path <- list.files(path = "data6", full.names = TRUE) %>%
  list.files(path = ., full.names = TRUE)

UCEC.mRNA <- readTCGA(path, dataType = "mRNA")

#####
##### miRNASeq
#####

# Directory in which untarred data will be stored
dir.create('data7')

# Download BRCA miRNASeq data and store it in data7 folder
# Remember that miRNASeq data are produced by two machines:
# Illumina Genome Analyzer and Illumina HiSeq 2000 machines
cancerType <- "BRCA"
downloadTCGA(cancerTypes = cancerType,
             dataSet = "Merge_mirnaseq__illuminaga_mirnaseq__bcgsc_ca__Level_3__miR_gene_expression__data.Level_3
             destDir = "data7")

downloadTCGA(cancerTypes = cancerType,
             dataSet = "Merge_mirnaseq__illuminahiseq_mirnaseq__bcgsc_ca__Level_3__miR_gene_expression__data.Level
             destDir = "data7")

# Shorten path of subdirectory with BRCA miRNASeq data
list.files(path = "data7", full.names = TRUE) %>%
  sapply(function(path){
    if (grepl(pattern = "illuminaga", path)){
      file.rename(from = grep(pattern = "illuminaga", path, value = TRUE),
                  to = file.path("data7",paste0(cancerType, ".miRNASeq.illuminaga")))
    } else if (grepl(pattern = "illuminahiseq", path)){
      file.rename(from = grep(pattern = "illuminahiseq", path, value = TRUE),

```

```

        to = file.path("data7",paste0(cancerType, ".miRNASeq.illuminahisec"))
    }
})

# Remove manifest.txt file
list.files(path = "data7", full.names = TRUE) %>%
  list.files(path = ., full.names = TRUE) %>%
  grep("MANIFEST.txt", x = ., value = TRUE) %>%
  file.remove()

# Read BRCA miRNASeq data
path <- list.files(path = "data7", full.names = TRUE) %>%
  list.files(path = ., full.names = TRUE)
path_illumina <- grep("illumina", path, fixed = TRUE, value = TRUE)
path_illuminahisec <- grep("illuminahisec", path, fixed = TRUE, value = TRUE)

BRCA.miRNASeq.illumina <- readTCGA(path_illumina, dataType = "miRNASeq")
BRCA.miRNASeq.illuminahisec <- readTCGA(path_illuminahisec, dataType = "miRNASeq")

BRCA.miRNASeq.illumina <- cbind(machine = "Illumina Genome Analyzer", BRCA.miRNASeq.illumina)
BRCA.miRNASeq.illuminahisec <- cbind(machine = "Illumina HiSeq 2000", BRCA.miRNASeq.illuminahisec)

BRCA.miRNASeq <- rbind(BRCA.miRNASeq.illumina, BRCA.miRNASeq.illuminahisec)

#####
#### isoforms
#####

# Directory in which untarred data will be stored
dir.create('data8')

# Download ACC isoforms data and store it in data8 folder
cancerType = "ACC"
downloadTCGA(cancerTypes = cancerType,
             dataSet = "Merge_rnaseqv2__illuminahisec_rnaseqv2__unc_edu__Level_3__RSEM_isoforms_normalized__data.",
             destDir = "data8")

# Shorten path of subdirectory with ACC isoforms data
list.files(path = "data8", full.names = TRUE) %>%
  file.rename(from = ., to = file.path("data8",paste0(cancerType, ".isoforms")))

# Remove manifest.txt file
list.files(path = "data8", full.names = TRUE) %>%
  list.files(path = ., full.names = TRUE) %>%
  grep("MANIFEST.txt", x = ., value = TRUE) %>%
  file.remove()

# Read ACC isoforms data
path <- list.files(path = "data8", full.names = TRUE) %>%
  list.files(path = ., full.names = TRUE)

ACC.isoforms <- readTCGA(path, dataType = "isoforms")

```

## End(Not run)

# Index

`as.data.frame`, 8

`checkTCGA`, 2, 3, 5–7, 9

`clinical`, 4

`datasetsTCGA`, 2, 3, 4, 6, 7, 9

`downloadTCGA`, 2, 3, 5, 6, 7–9

`infoTCGA`, 2, 3, 5, 6, 7, 9

`mutations`, 4

`readTCGA`, 2, 3, 5–7, 8

`rnaseq`, 4

`RTCGA (RTCGA-package)`, 2

`RTCGA-package`, 2