



Databases and ontologies

# RTCGA - The Family of R Packages Integrating Data from The Cancer Genome Atlas Study

Marcin Kosiński<sup>1,2,\*</sup>, Witold Chodor<sup>2</sup> and Przemysław Biecek<sup>1,2</sup>

<sup>1</sup>Faculty of Mathematics and Information Science, Warsaw University of Technology, 00-662 Warsaw, Poland and

<sup>2</sup>Faculty of Mathematics, Informatics, and Mechanics, University of Warsaw, City, 02-097 Warsaw, Poland.

\*To whom correspondence should be addressed.

Associate Editor: XXXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

## Abstract

**Summary:** In this article we present a family of **R** packages called **RTCGA** that facilitate access to data from the TCGA project. The Cancer Genome Atlas Project (TCGA) is a coordinated effort to accelerate our understanding of the molecular basis of cancer. It is a source of curated multi-platform data, including RNA-seq, DNA-seq, DNA Methylation, together with clinical data for over 11 thousand patients and 33 cancer types. This rich source of data is accessible in raw format from TCGA Data Portal. **RTCGA** packages facilitate access to this dataset, merge patients characteristics from different platforms and support typical statistical analyses. These packages will be useful for at least three audiences: biostatisticians that work with cancer data; researchers that are engaged with large scale algorithms; lecturers that are presenting data analysis methods on real data problems.

**Availability:** **RTCGA** family of **R** packages is freely available at GitHub <http://rtcga.github.io/RTCGA/> and from the Bioconductor project at <http://bioconductor.org/packages/RTCGA/>.

**Contact:** m.p.kosinski@gmail.com

## Motivation

The Cancer Genome Atlas Data Portal provides a platform for researchers to search, download, and analyze data sets generated by TCGA Project. It contains clinical information, genomic characterization data, and high level sequence analysis of the tumor genomes. Files are available through Firehose Broad GDAC portal. One can select cancer type (cohort) and data type (e.g. clinical, RNA expression, methylation) and download a 'tar.gz' file with compressed data. While working with many cancer types we may find this approach burdensome. For reasons listed further we prepared selected datasets from the TCGA project in an easy to process way and embedded them in separate R data packages.

- If one requires to download datasets containing e.g. information about genes' expressions for all available cohorts types, then one would have to go through the click-to-download process many times. This is inconvenient and time-consuming.
- Some datasets (e.g. clinical) are not in a standard tidy data format, which is: one row for one observation and one column for one variable. They are transposed. That becomes more onerous when investigating many clinical datasets at once.

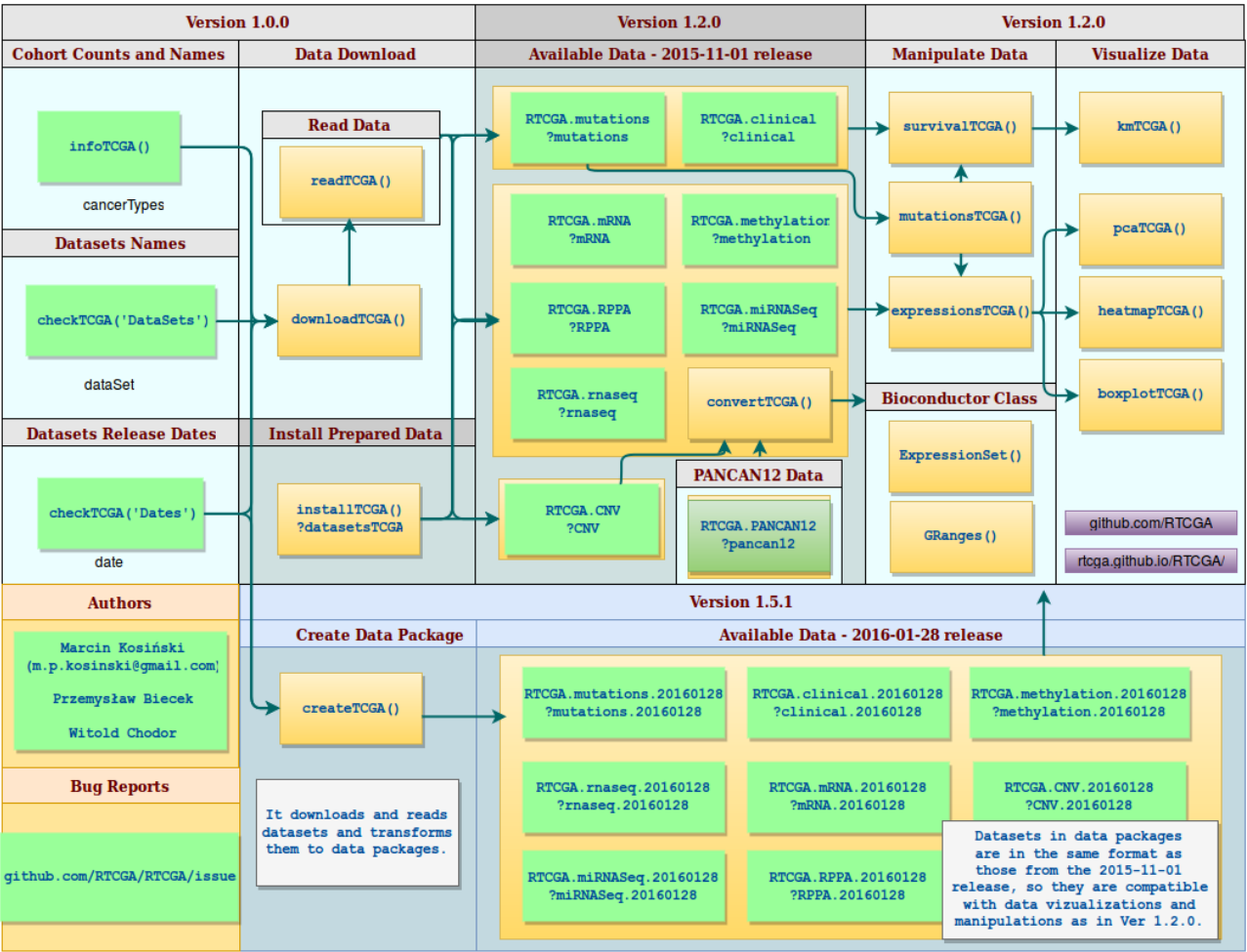
- Datasets containing information on some data types (e.g. gene's mutations) are not in one easy-to-handle file. Every patient has its own file, which may be an impassable barrier for many potential researchers.
- Data governance for many datasets for various cohorts saved in different folders with strange (default after untarring) names may be exhausting and uncomfortable for researchers that are not very skilled in data management or data processing.

## Architecture

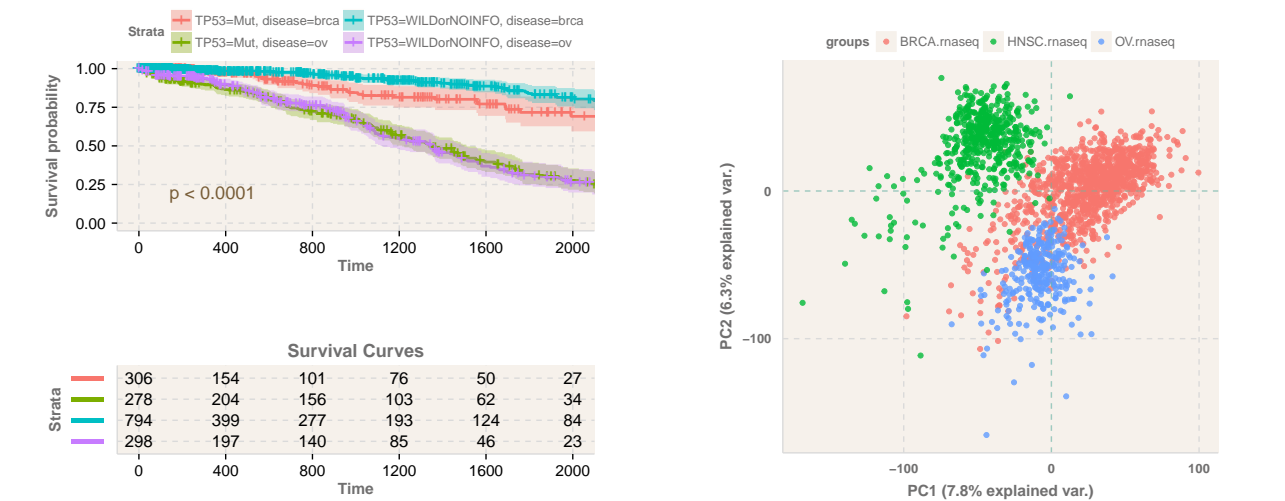
The general architecture of **RTCGA** is presented in the Figure 1a. The software package contains functions that facilitate download of data from particular date, cohort and platform. Other functions benefits data processing, analysis and visualizations. In figures 1b-1c we present example analyses and plots that cover Kaplan-Meier estimates of survival curves and Principal Component Analysis. Detailed instructions on how to apply these analyses (and many other) and on how to download data are presented at the project's website.

## References

Broad Institute of MIT and Harvard (2014). *Broad Institute TCGA Genome Data Analysis Center: Firehose stddata 2015 06 01 run*. DOI:10.7908/C1251HBG.



(a) The workflow of **RTCGA** family of software and data R packages. **RTCGA** consists of auxiliary functions: `infoTCGA()` and `checkTCGA` that enable to find out metadata about datasets and dates of their release provided by TCGA. Data download section consists of `downloadTCGA()` function that allows to download every dataset from TCGA study and `readTCGA()` function that enables to read (most popular) data into the tidy format. It is possible to use `installTCGA()` function to install, prepared in tidy format, most popular data types from TCGA that are included in data packages. Datasets in **RTCGA** data packages can be converted to Bioconductor format (`ExpressionSet`, `GRanges`) with `convertTCGA()` function. Results of few functions designed to manipulate and visualize **RTCGA** data are presented on subfigures (b)-(c). In version 1.5.1 the `createTCGA()` function was added, which allows to create TCGA R data packages.



(b) Kaplan-Meier estimates of survival curves and risk set table for **Breast invasive carcinoma (BRCA)** and **Ovarian serous cystadenocarcinoma (OV)** cohorts divided on **TP53** mutations. The effect of `survivalTCGA()`, `mutationsTCGA()` and `kmTCGA()` functions. (c) Principal Component Analysis performed for genes expressions (RNA-Seq) for **Breast invasive carcinoma (BRCA)**, **Head and Neck squamous cell carcinoma (HNSC)** and **Ovarian serous cystadenocarcinoma (OV)** cohorts. The effect of `expressionsTCGA()` and `pcaTCGA()`.

Fig. 1: The workflow of **RTCGA** family of software and data R packages and effects of functions designed to manipulate and visualize **RTCGA** data.