

UCSCXenaShiny: an R package for exploring and analyzing UCSC Xena public datasets in web browser

Shixiang Wang^{*,a,b,c,d}, Yi Xiong^{*,a,e,f}, Kai Gu^{*,a,g}, Longfei Zhao^{a,h}, Yin Li^{a,j},
Fei Zhao^{a,d,i}, Xuejun Li^f, Xue-Song Liu^{**,b}

^aOpenbioX Community, China

^bSchool of Life Science and Technology, ShanghaiTech University, China

^cShanghai Institute of Biochemistry and Cell Biology, Chinese Academy of Sciences, China

^dUniversity of Chinese Academy of Sciences, China

^eXiangya School of Medicine, Central South University, China

^fDepartment of Neurosurgery, Xiangya Hospital, Central South University, China

^gRoche Diagnostics (Shanghai) Limited, China

^hZhengzhou University, China

ⁱCAS Center for Excellence in Molecular Plant Sciences, China

^jDepartment of Thoracic Surgery, Zhongshan hospital, Fudan university, China

Abstract

Motivation: UCSC Xena platform provides huge amounts of processed cancer omics data from big public projects like TCGA or individual reserach groups for enabling unprecedented research opportunities. In 2019, we developed UCSCXenaTools, an R package for retrieval of UCSC Xena data. However, an easier dataset exploration and analysis tool is still lack, especially for researchers without programming experience.

Results: We develop UCSCXenaShiny, an R Shiny package to quickly explore, download all datasets from UCSC Xena data hubs. In addition, a module based analysis framework is constructed to analyze and visualize data.

Availability: <https://github.com/openbioX/UCSCXenaShiny> or <https://cran.r-project.org/package=UCSCXenaShiny>.

Key words: UCSC-Xena cancer-genomics TCGA

Introduction

Over the past decade, programs including TCGA (Weinstein et al., 2013), ICGC (Zhang et al., 2011), PCAWG (The et al., 2020), GTEx (Consortium and others, 2015), CCLE (Barretina et al., 2012) and etc. have generated large amounts of molecular data characterizing the landscape of more than ten

*Equal contribution

**Corresponding Author

Email addresses: wangshx@shanghaitech.edu.cn (Shixiang Wang),
xiongyi123@csu.edu.cn (Yi Xiong), gukai1212@163.com (Kai Gu), longfei8533@live.cn
(Longfei Zhao), yinli@openbioX.org (Yin Li), zhaofei@openbioX.org (Fei Zhao),
lxjneuro@csu.edu.cn (Xuejun Li), liuxs@shanghaitech.edu.cn (Xue-Song Liu)

Preprint submitted to Preprints

July 9, 2020

thousands of tumors from genomic, epigenetic and proteomic aspects. The data have been preprocessed and stored at data hubs of UCSC Xena platform along with many public cancer datasets from individual research groups, providing unprecedented opportunities for either simple or systematic exploration of cancer behaviors and mechanisms at multiple molecular layers in individual cancer type or across cancer types (Goldman et al., 2019).

In 2019, we developed UCSCXenaTools, an open-source R package for retrieving and assembling public UCSC Xena data (Wang and Liu, 2019). UCSCXenaTools was developed to communicate with UCSC Xena data hubs for downloading datasets or dataset subsets, querying metadata of data hub, cohort or dataset. Despite UCSC Xena platform itself allows users to explore and analyze data, it is hard for researchers to quickly explore all available datasets, locate what they need in their research and download useful datasets. Besides, the analysis features provided by UCSC Xena platform mainly focus on individual cohort data, thus lack of full-feature functionality.

To this end, we develop an open-source R Shiny package UCSCXenaShiny for cancer community to allow researchers to explore and analyze datasets from UCSC Xena data hubs in web browser. In addition, an extensible module based analysis framework is constructed to analyze data. Currently, several modules providing single-gene expression analysis and visualization are implemented.

Materials and methods

Dataset exploration

UCSCXenaShiny opens a web page in user's browser to provide service. The page "Repository" is used to explore all available UCSC Xena datasets. Users can find desired datasets by either defined buttons or searching in dataset table. Once one or several datasets selected, users can query their metadata or download them (Fig.1 and Fig.3). To improve the performance of downloading large datasets, we provide a button to download a Shell script containing 'wget' commands which can run in Unix-like system.

Module and pipeline

For now, several modules targeting at single-gene expression analysis are available at page "Module" (Fig.2 and Fig.4), a pipeline based on them is available at page "Pipeline" (Fig.2). The usage is quite easy, users just need to type the gene symbol name and all procedures will be properly done by UCSCXenaShiny, including downloading data from UCSC Xena data hubs, cleaning data, analyzing data and visualizing the result. We are happy to accept new feature requests and they can be discussed at <https://github.com/openbioX/UCSCXenaShiny/issues>.

Results

Package structure

The structure and workflow of UCSCXenaShiny is described in Fig.1. Currently, the core components of this package are page "Repository" and page

“Module”. Page “Repository” allows researchers to explore and download datasets. Table 1 summaries the cohort and dataset number available at different UCSC Xena data hubs. There are total 1639 datasets and TCGA project is the major contributor. The development of UCSCXenaShiny is based on R Shiny platform (<https://shiny.rstudio.com/>), the overview of its graphic interface is shown in Fig.2.

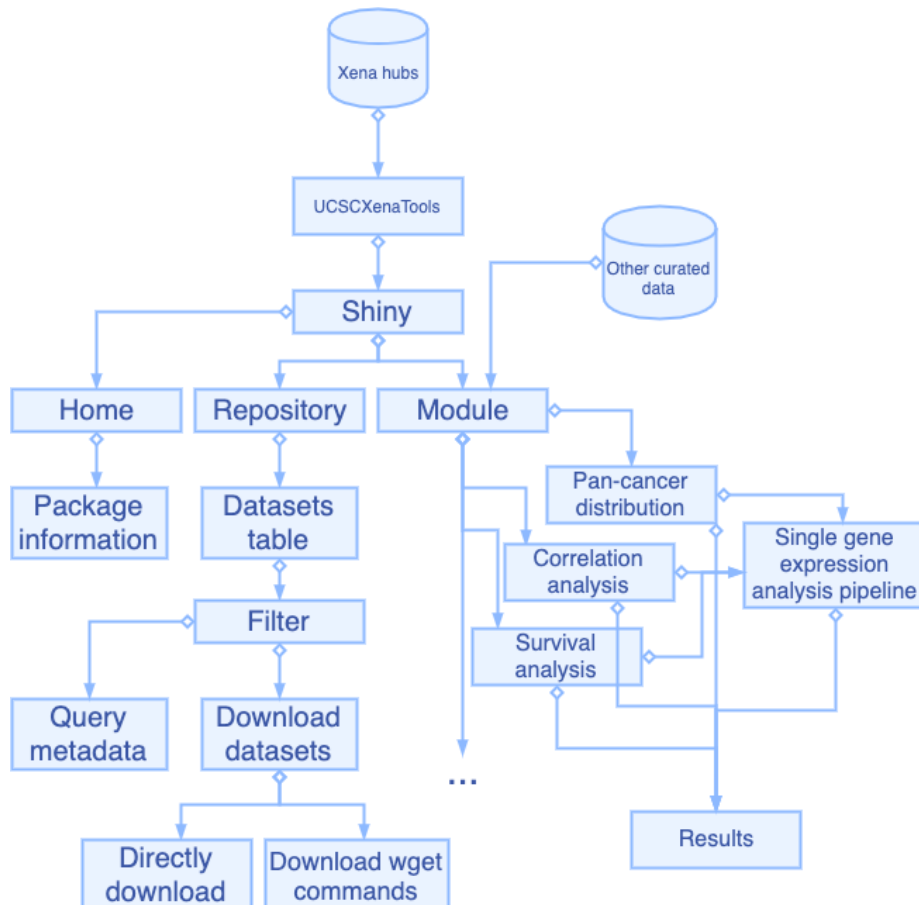


Figure 1: Package architecture and functional flowchart of UCSCXenaShiny.

Feature 1: dataset exploration and download

UCSCXenaShiny allows users to explore UCSC Xena datasets quickly and easily (Fig.3). A table storing all datasets is shown in Fig.3A, users can filter datasets by either typing some key words in search bar or selecting data hubs or data types. Once desired datasets are selected in the table, users can click the button on the bottom to check metadata of datasets or download datasets (Fig.3B-D).

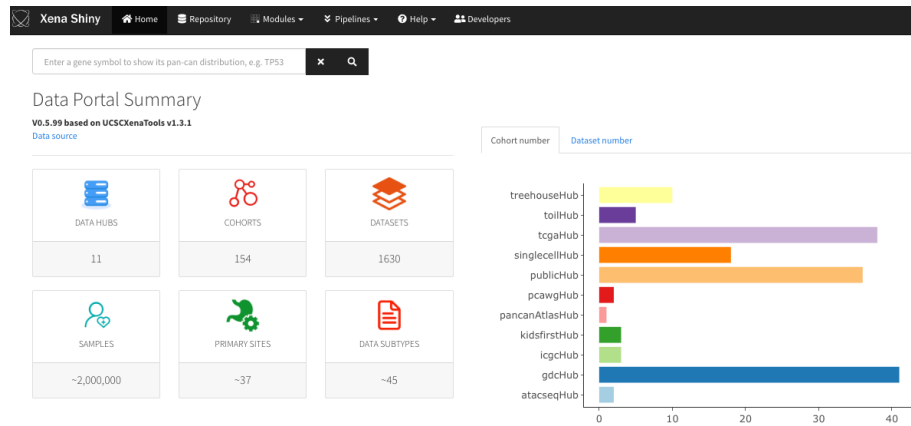


Figure 2: UCSCXenaShiny graphical interface overview.

	Data hub	Cohorts	Datasets	URL
1	tcgaHub	38	715	https://tcga.xenahubs.net
2	gdcHub	41	528	https://gdc.xenahubs.net
3	publicHub	36	109	https://ucscpublic.xenahubs.net
4	pcawgHub	2	53	https://pcawg.xenahubs.net
5	toilHub	5	50	https://toil.xenahubs.net
6	singlecellHub	18	54	https://singlecellnew.xenahubs.net
7	icgcHub	3	23	https://icgc.xenahubs.net
8	pancanAtlasHub	1	22	https://pancanatlas.xenahubs.net
9	treehouseHub	10	26	https://xena.treehouse.gi.ucsc.edu
10	atacseqHub	2	9	https://atacseq.xenahubs.net
11	kidsfirstHub	3	50	https://kidsfirst.xenahubs.net

Table 1: Summary of UCSC Xena data hubs.

Feature 2: Single-gene expression analysis

UCSCXenaShiny provides modules implementing basic analysis functionality and modules can be go further assembled as analysis pipeline (Fig.1). For example, we constructed a few modules to analyze and visualize the single gene expression, including its pan-cancer distribution with violin plot or anatomy heatmap (Maag, 2018), and survival effects (Terry M. Therneau and Patricia M. Grambsch, 2000) under different expression cutoff. We combined some of them and built single-gene expression analysis pipeline so researchers can get as much information as possible in one click for a same task view. An example for gene *TP53* is given in Fig.4.

Acknowledgements

We thank projects TCGA, GTEx, ICGC, CCLE, PCAWG, etc., and individual research groups for making cancer genomics data public. We thank UCSC Xena platform for providing data processing, integration and download.

References

- Barretina, J., Caponigro, G., Stransky, N., Venkatesan, K., Margolin, A.A., Kim, S., Wilson, C.J., Lehár, J., Kryukov, G.V., Sonkin, D., others, 2012. The cancer cell line encyclopedia enables predictive modelling of anticancer drug sensitivity. *Nature* 483, 603.
- Consortium, G., others, 2015. The genotype-tissue expression (gtex) pilot analysis: Multitissue gene regulation in humans. *Science* 348, 648–660.
- Goldman, M., Craft, B., Hastie, M., Repečka, K., Kamath, A., McDade, F., Rogers, D., Brooks, A.N., Zhu, J., Haussler, D., 2019. The ucsc xena platform for cancer genomics data visualization and interpretation. *BioRxiv* 326470.
- Maag, J., 2018. Gganatogram: An r package for modular visualisation of anagrams and tissues based on ggplot2. *f1000research*.
- Terry M. Therneau, Patricia M. Grambsch, 2000. *Modeling survival data: Extending the Cox model*. Springer, New York.
- The, I., Whole, T.P.-C.A. of, Consortium, G., others, 2020. Pan-cancer analysis of whole genomes. *Nature* 578, 82.
- Wang, S., Liu, X., 2019. The ucscxenatools r package: A toolkit for accessing genomics data from ucsc xena platform, from cancer multi-omics to single-cell rna-seq. *Journal of Open Source Software* 4, 1627.
- Weinstein, J.N., Collisson, E.A., Mills, G.B., Shaw, K.R.M., Ozenberger, B.A., Ellrott, K., Shmulevich, I., Sander, C., Stuart, J.M., Network, C.G.A.R., others, 2013. The cancer genome atlas pan-cancer analysis project. *Nature genetics* 45, 1113.
- Zhang, J., Baran, J., Cros, A., Guberman, J.M., Haider, S., Hsu, J., Liang, Y., Rivkin, E., Wang, J., Whitty, B., others, 2011. International cancer genome consortium data portal—a one-stop shop for cancer genomics data. *Database* 2011.

A

Dataset Filters

Active Data Hub:

- UCSF Public
- TCGA
- GDC
- ICGC
- Pan-Cancer Atlas
- TOIL
- Treehouse
- PCAWG
- ATAC-seq
- Single Cell

Cohort Name:

e.g. breast (separator &)

Data Type:

- Phenotype
- Feature by sample matrix
- Genomic segments
- Mutations

Data Subtype:

e.g. gene expression (separator &)

How to use repository

Show 10 entries

Dataset	Hub	Cohort	Samples	Subtype	Label	Unit
1 TCGA-BLCA.cmv.tsv	gdcHub	GDC TCGA Bladder Cancer (BLCA)	415	copy number	Copy Number Segment	log2(copy number/2)
2 TCGA-BLCA.gdc_phenotype.tsv	gdcHub	GDC TCGA Bladder Cancer (BLCA)	454	phenotype	Phenotype	
3 TCGA-BLCA.gistic.tsv	gdcHub	GDC TCGA Bladder Cancer (BLCA)	413	copy number (gene-level)	GISTIC - focal score by gene	Gistic2 copy number
4 TCGA-BLCA.htseq_counts.tsv	gdcHub	GDC TCGA Bladder Cancer (BLCA)	430	gene expression RNAseq	HTSeq - Counts	log2(count+1)
5 TCGA-BLCA.htseq_fpkm.tsv	gdcHub	GDC TCGA Bladder Cancer (BLCA)	430	gene expression RNAseq	HTSeq - FPKM	log2(fpkm+1)
6 TCGA-BLCA.htseq_fpkm_uq.tsv	gdcHub	GDC TCGA Bladder Cancer (BLCA)	430	gene expression RNAseq	HTSeq - FPKM UQ	log2(fpkm+1)
7 TCGA-BLCA.masked_cmv.tsv	gdcHub	GDC TCGA Bladder Cancer (BLCA)	415	copy number	Masked Copy Number Segment	log2(copy number/2)
8 TCGA-BLCA.methylation450.tsv	gdcHub	GDC TCGA Bladder Cancer (BLCA)	437	DNA methylation	Illumina Human Methylation 450	beta value
9 TCGA-BLCA.mimaz.tsv	gdcHub	GDC TCGA Bladder Cancer (BLCA)	432	stem loop expression	miRNA Expression Quantification	log2(RPM+1)
10 TCGA-BLCA.mutuse_srv.tsv	gdcHub	GDC TCGA Bladder Cancer (BLCA)	411	somatic mutation (SNPs and small INDELS)	MuSE Variant Aggregation and Masking	

Showing 1 to 10 of 528 entries

Previous 1 2 3 4 5 ... 53 Next

Show Metadata Request Data

B

Submitted datasets:

Index	URL	Dataset	Unit
1	https://gdc.xenahubs.net	TCGA-BLCA.GDC_phenotype.tsv	https://gdc.xenahubs.net/download/TCGA-BLCA.GDC_phenotype.tsv.gz
2	https://gdc.xenahubs.net	TCGA-BLCA.htseq_counts.tsv	https://gdc.xenahubs.net/download/TCGA-BLCA.htseq_counts.tsv.gz

Download data directly Batch download in terminal

C

Please select a folder

Create new folder Sort content

home

Directories

- home
- anaconda3
- Applications
- biodata
- biosoft
- Desktop
- Documents
- Downloads
- fsdownload
- go
- icmp4j
- igv
- Library
- MessageBoard
- Movies
- Music
- Nutstore Files
- OneDrive - shanghaiitech.edu.cn
- Pictures
- Public
- quicklisp
- R-dev

Content

No folder selected

Cancel Select

D

2020-07-06-commands.sh — Downloads

```

1 #!/usr/bin/env bash
2 #Baller run bash 2020-07-06-commands.sh in your terminal under a desired
3 directory
4 wget -c https://gdc.xenahubs.net/download/TCGA-BLCA.GDC_phenotype.tsv.gz
5 wget -c https://gdc.xenahubs.net/download/TCGA-BLCA.htseq_counts.tsv.gz

```

Line: 1 Shell Script (Bash) Tab Size: 4

Figure 3: Dataset search and download.

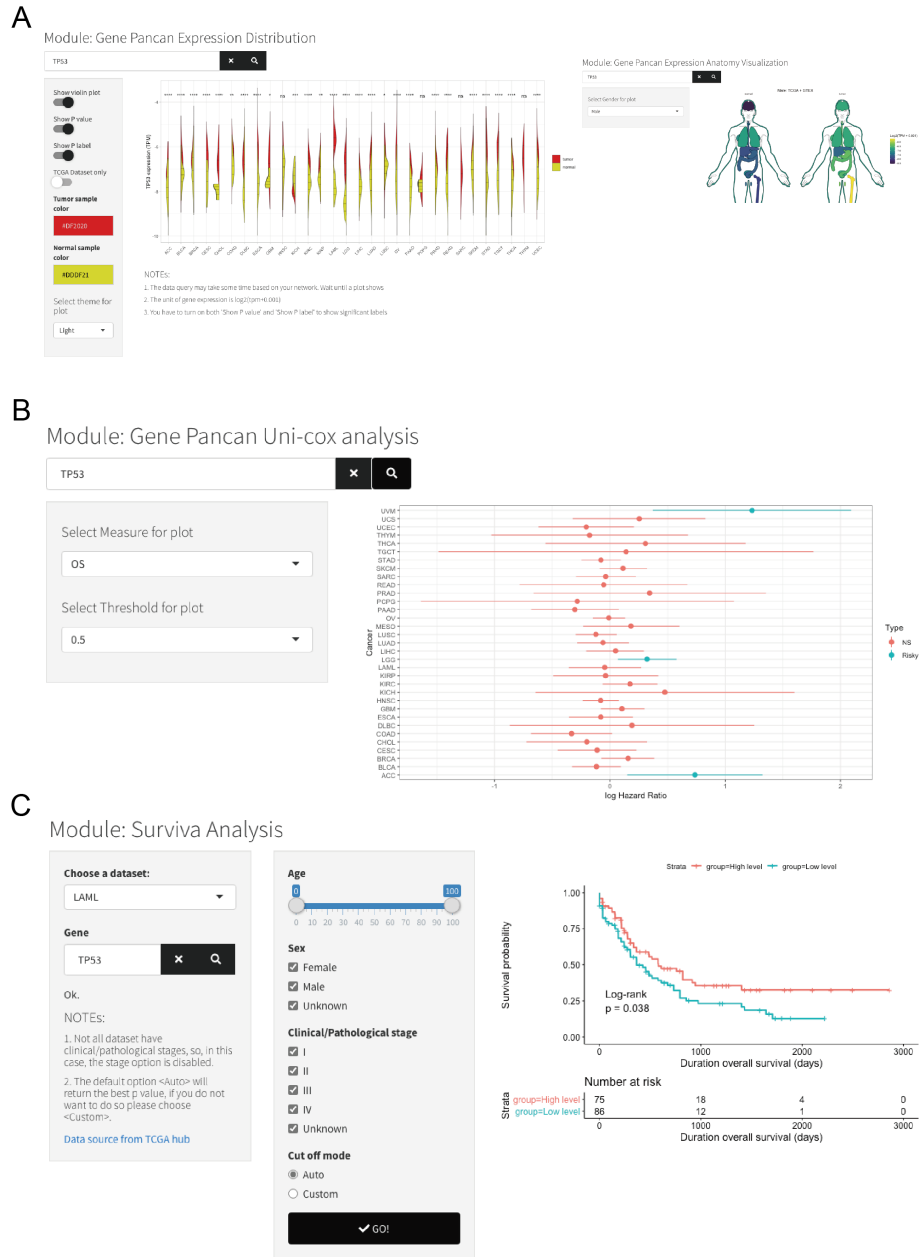


Figure 4: Current available analysis modules provided UCSCXenaShiny.