# Tutorial

Version 2016_04_17

# Introduction

A simple and elegant way to explore cancer data.

Backed by a powerful computational infrastructure, application programming interface (API), graphical tools and online reports.

Sitting above one of the deepest and most integratively-characterized **_open_** cancer datasets in the world.

With over 80K sample aliquots from 11,000+ cancer patients, spanning 38 unique disease cohorts.

**FIREBROWSE**

Data and analyses utilized at numerous academic, research, and commercial sites around the world.

Example: cBio@MSKCC

TCGA data & analyses in cBioPortal—expression, mutation, copy number, signficance analyses, and more—are loaded directly from Firehose.

**FIREBROWSE**

Search analysis results

View Expression Profile | Enter gene name | Enter cohort abbrev | View Analysis Profile

**SELECT COHORT** ▼

- Clinical Analyses
- CopyNumber Analyses
- Correlations Analyses
- miR Analyses
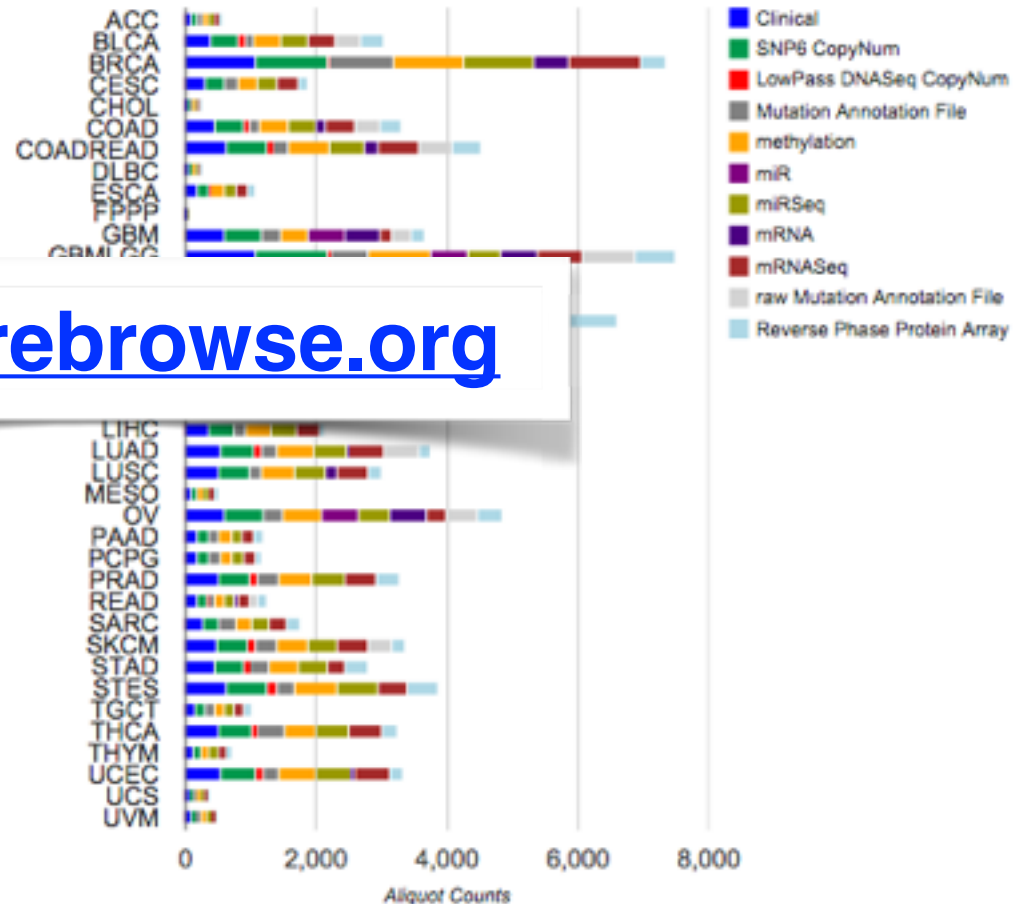- miRseq Analyses
- mRNA Analyses
- mRNAseq Analyses
- Mutation Analyses
- Pathway Analyses
- RPPA Analyses

TCGA data version 2015_06_01

Legend:
- Clinical
- SNP6 CopyNum
- LowPass DNASeq CopyNum
- Mutation Annotation File
- methylation
- miR
- miRSeq
- mRNA
- mRNASeq
- raw Mutation Annotation File
- Reverse Phase Protein Array

Cohorts: ACC, BLCA, BRCA, CESC, CHOL, COAD, COADREAD, DLBC, ESCA, FPPP, GBM, GBMLGG, ... LIHC, LUAD, LUSC, MESO, OV, PAAD, PCPG, PRAD, READ, SARC, SKCM, STAD, STES, TGCT, THCA, THYM, UCEC, UCS, UVM

X-axis: 0, 2,000, 4,000, 6,000, 8,000 — *Aliquot Counts*

## http://firebrowse.org

# Many 1000s of datasets per run
# Find your favorite in 2 clicks

**Choose Cohort**

**Then DataType**

**Click to download**

Thyroid carcinoma (THCA)

TCGA data version 2016_01_28 for THCA

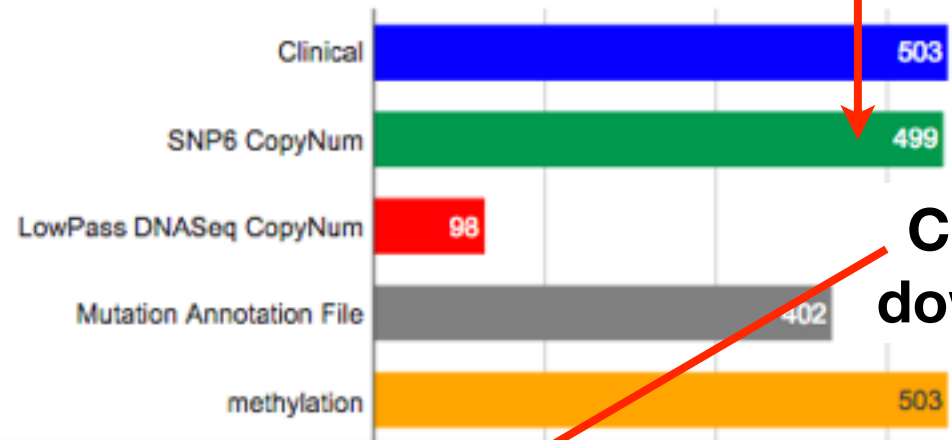| ■ Clinical Analyses |
| ■ CopyNumber Analyses |
| Correlations Analyses |
| ■ Methylation Analyses |
| ■ miRseq Analyses |
| ■ mRNA A |
| ■ mRNAse |
| ■ Mutation |
| Pathway |
| ■ RPPA An |

Clinical — 503
SNP6 CopyNum — 499
LowPass DNASeq CopyNum — 98
Mutation Annotation File — 402
methylation — 503

## THCA CopyNumber Archives

| Primary | Auxiliary | SDRF/Mage | Send To |

Files may also be downloaded here, or with firehose_get, or exported to GenomeSpace with the SendTo tab.

genome_wide_snp_6-segmented_scna_minus_germline_cnv_hg19  (MD5)
genome_wide_snp_6-segmented_scna_hg19  (MD5)
genome_wide_snp_6-segmented_scna_minus_germline_cnv_hg18  (MD5)
genome_wide_snp_6-segmented_scna_hg18  (MD5)

Downloading data constitutes agreement to TCGA data usage policy

600

# Or easily send to GenomeSpace for more analysis



## THCA CopyNumber Archives

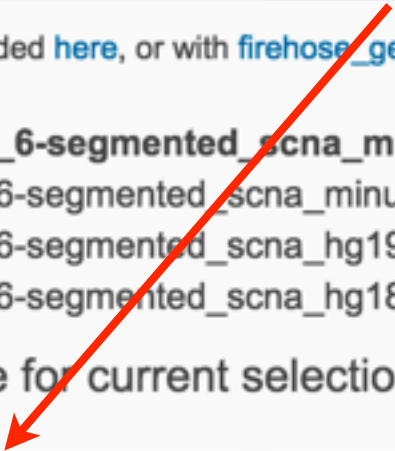| Primary | Auxiliary | SDRF/Mage | **Send To** |

Files may also be downloaded here, or with firehose_get, or exported to GenomeSpace with the SendTo tab.

☑ **genome_wide_snp_6-segmented_scna_minus_germline_cnv_hg19 [851.18 KB]**
☐ genome_wide_snp_6-segmented_scna_minus_germline_cnv_hg18 [852.54 KB]
☐ genome_wide_snp_6-segmented_scna_hg19 [6.46 MB]
☐ genome_wide_snp_6-segmented_scna_hg18 [6.47 MB]

Cumulative file size for current selections: 851.18 KB

| ☁️ GENOMESPACE Upload | Clear Selections |

Downloading data constitutes agreement to TCGA data usage policy

# GENOMESPACE

Frictionless connection of bioinformatics tools

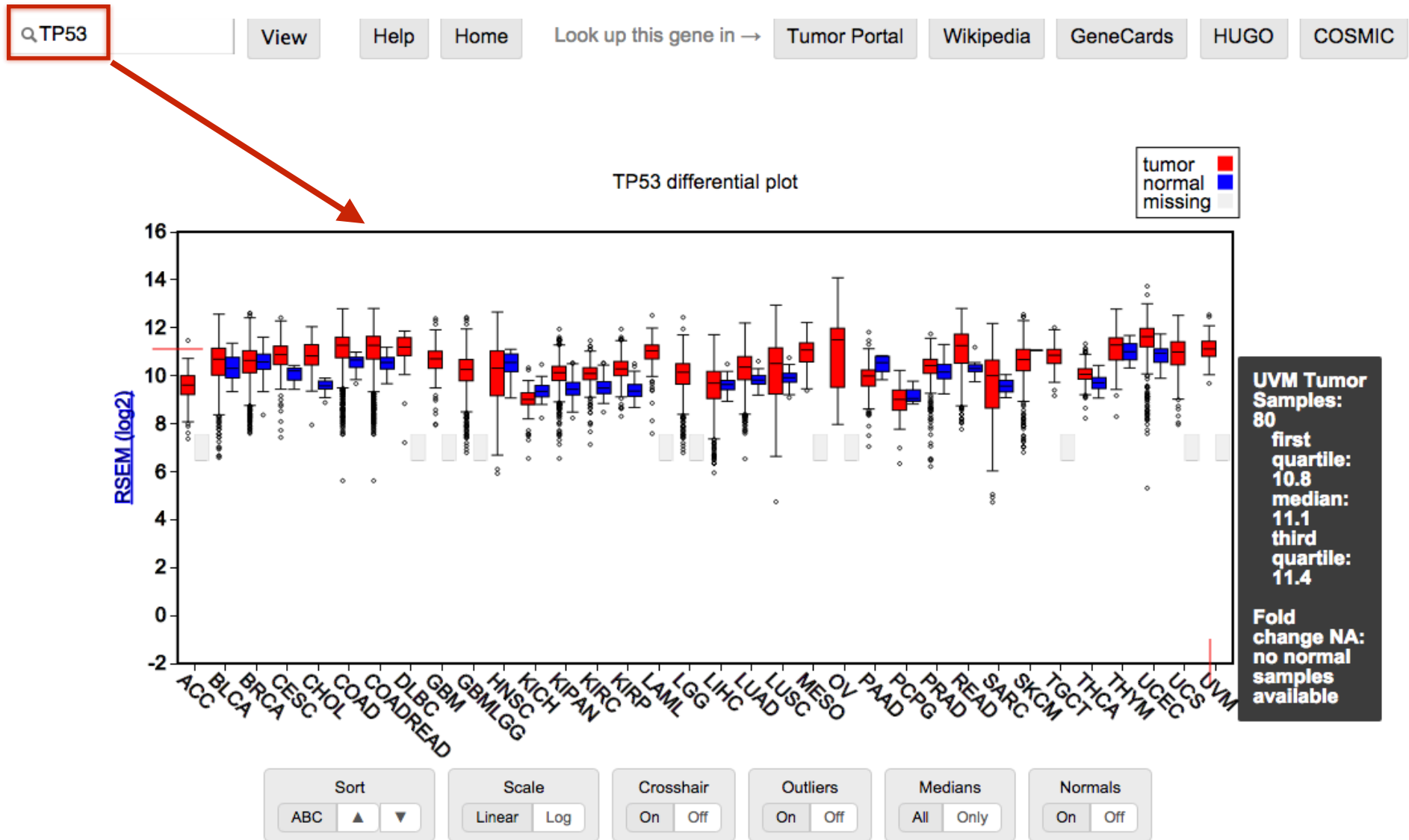| Register | User Login |

## Or download everything with 1 command

```
linux% firehose_get analyses latest
```

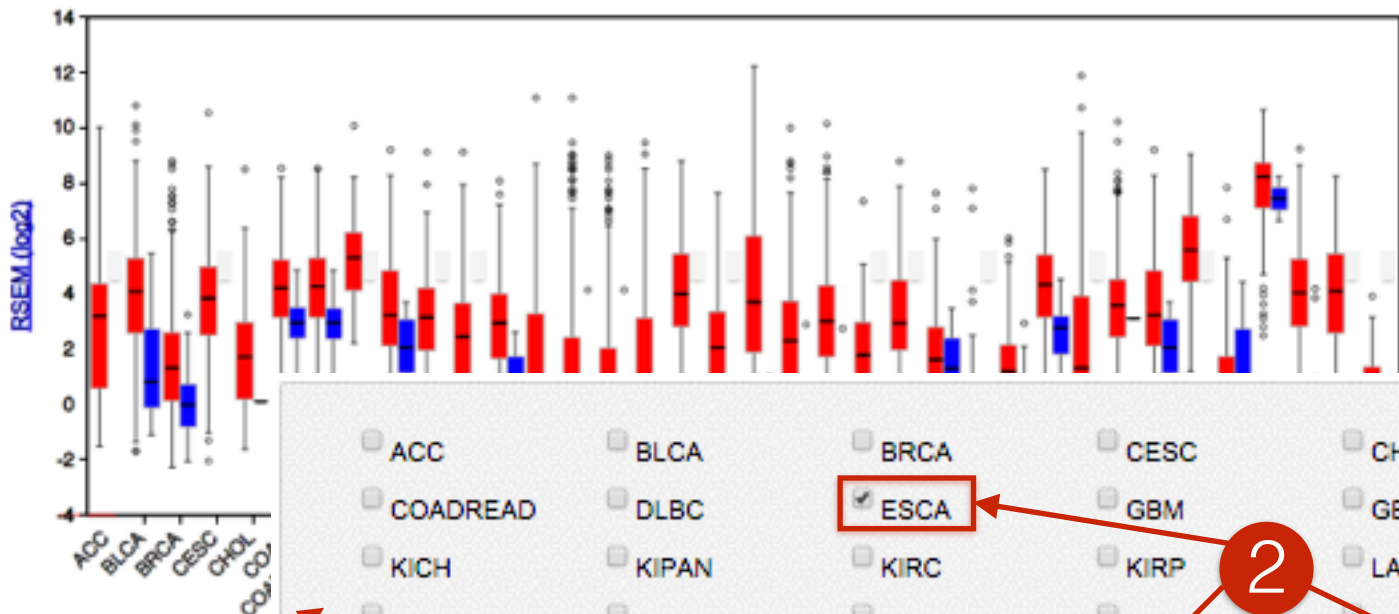**Simple 20K bash script, just 1 moving part**

[Download Here](#)

# Graphical Tools

# viewGene: expression level browser



Quickly inspect mRNASeq expression levels for a selected gene

TERT differential plot

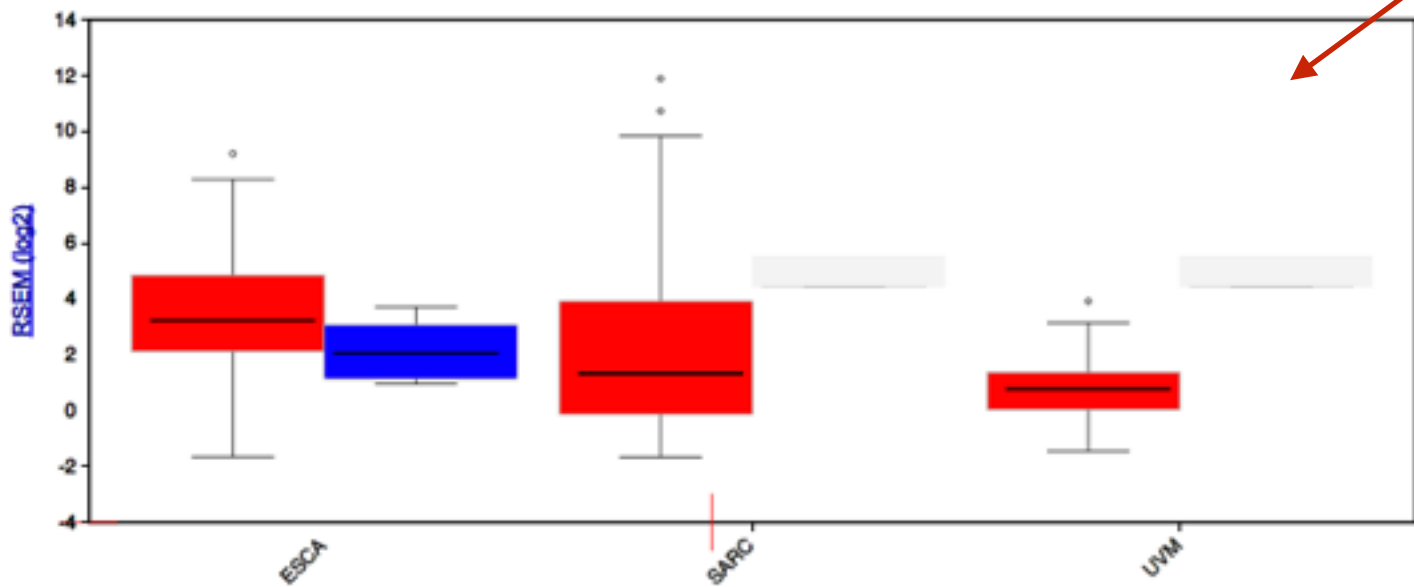View expression levels across all cohorts, or arbitrary subsets.

# CoMut: mutation co-occurence plots



Introduced in 2011 (Stransky et al, Science, 2011), CoMut figures have become common in TCGA research. Within a single graphic they provide a _comprehensive analysis profile_, enabling the reader to quickly infer relationships between co-occurring results across multiple data modalities, across common X axis of sample IDs.

But in journals, figures are static and can be small and hard to read

And cannot be explored in real time

And reproducing them or investigating their implications can require substantial time for data retrieval, preparation and analysis

Now we've re-sorted by CNMF copy-number clustering, ***and dragged it from bottom of figure to top,*** just above mutation panel

Making it further apparent that the copy-number landscape differs as IDH1/2, TP53, and ATRX mutations diminish

Also shows apparent involvement with EGFR and PTEN.

# iCoMut compressess an entire Firehose run into a single, interactive & reproducible figure



**Turning this …**

Firehose analysis workflow

Run on 38 TCGA cohorts
> 100 tasks per run
~10 datatypes

Distills 50 TB of input data
into 10GB of results (5000x)

View Expression Profile | Enter gene name 🔍 | UVM 🔍 | View Analysis Profile

… into this

**firebrowse.org/iCoMut/?cohort=UVM**

# By default, patients are sorted by histology and gene mutation

# Sort status of samples is reported in the info box

**CESC**                                                    ×

**# Samples: 194 patients**

**Samples are sorted by ...**
histology, PIK3CA, EP300, FBXW7, HLA.A,
ARID1A, PTEN, MTOR*, FAT2, NFE2L2, NHS,
KRAS, ERBB2*, MED1, ERBB3, HLA.B, TP53*,
ZNF750, TRIM9, MAPK1, RB1*

Close

General
Help

iCoMut Beta for FireBrowse

CESC - Cervical squamous cell carcinoma and endocervical adenocarcinoma ▾

🏠 ℹ️ 🔄 🔳 ✖️ ❓

Mutation Rate        Mutation Rate    60
  • synonymous
  • non synonymous                   40

                                     20

                                      0

Clinical Age                   age
Clinical Vital Status     vital status
Clinical Gender             gender
➤ Clinical Histology        histology
Clinical Ethnicity            race

➤ Gene Mutation

Show Grids    Hide Grids    Multi-Key Sort

  • NA
  • Nonsense
  • Frameshift
  • Splice Site
  • Missense
  • Other Non Syn
  • In-frame INDEL
  • Syn
    No Mutation

PIK3CA
EP300
FBXW7
HLA.A
ARID1A
PTEN
MTOR*
FAT2
NFE2L2
NHS
KRAS
ERBB2*
MED1
ERBB3
HLA.B
TP53*
ZNF750
TRIM9
MAPK1
RB1*

40   20   0
# Mutations

# Push-Button Publication Figure Reproducibility



**Download Figure**

SVG

---

# iCoMut results for ACC - Adrenocortical carcinoma    *Disease Type*

Generated on Wed Nov 04 2015 1:34:46 PM with integ-2015_10_29_c | 9851a395fb32    *Date & Software Version*

View on Firebrowse at:

http://fbdev/iCoMut/?cohort=acc

**URL to regenerate (will reflect all interactive manipulations to figure)**

# Many more graphical controls ...



**Example: locate patient/sample of interest**

**Collaboratively explore questions in realtime on telecons:
in what expression cluster does patient X fall?**

**Without database lookup or scripting, etc**

# Advanced Search

×

Include these samples:
OR-A5K5

Exclude these samples:
C5-A0TN

## Panel Functions

| | Row | OP | Value |
|---|---|---|---|
| -- Select a Panel -- | ✓ non_synonymous | ✓ > | |
| ✓ mutation_rate | synonymous | >= | |
| clinical_age | total | < | |
| clinical_vital_status | | <= | |
| clinical_gender | ● AND ○ OR | = | Search |
| clinical_histology | | != | |
| clinical_ethnicity | | | |
| gene_mutation | | | |
| focal_level_cn_gain | | | |
| focal_level_cn_loss | | | |
| mrnaseq_cnmf | | | |
| mrnaseq_chierarchical | | | |
| mirseq_cnmf | | | |
| mirseq_cHierarchical | | | |
| mirseq_mature_cnmf | | | |
| mirseq_mature_chierarchical | | | |
| cn_cnmf | | | |
| clus_methylation_cnmf | | | |
| rppa_cnmf_clusters | | | |
| rppa_chierarchical | | | |

Drag and drop the − or + icon to rearrange the panels

Rearranged panels

# iCoMut Beta for FireBrowse

ACC - Adrenocortical carcinoma ▾

🏠 ❶ ↻ ⤢ ↙ ❓

ⓘ ⊞ Mutation Rate

ⓘ ⊞ Clinical Age

ⓘ ⊞ Clinical Vital Status

ⓘ ⊞ Clinical Gender

ⓘ ⊞ Clinical Histology

ⓘ ⊞ Clinical Ethnicity

ⓘ ⊞ Gene Mutation

ⓘ ⊞ Focal Level CN Gain

ⓘ ⊞ Focal Level CN Loss

ⓘ ⊞ mRNAseq cNMF

ⓘ ⊞ mRNAseq cHierarchical

ⓘ ⊞ miRseq cNMF

ⓘ ⊞ miRseq cHierarchical

ⓘ ⊞ miRseq Mature cNMF

ⓘ ⊞ miRseq Mature cHierarchical

ⓘ ⊞ CN cNMF

ⓘ ⊞ Methylation cNMF

ⓘ ⊞ RPPA cNMF Clusters

ⓘ ⊞ RPPA cHierarchical

More features described in online help

**http://firebrowse.org/iCoMut/#icomutHelp**

# Programmatic Tools

# API-Powered :  25+ RESTful apis in 4 categories

**Analyses** : Fine grained retrieval of analysis pipeline results    Show/Hide | List Operations | Expand Operations | Raw

| GET | /Analyses/Mutation/MAF | Retrieve MutSig final analysis MAF. |
| GET | /Analyses/Mutation/SMG | Retrieve Significantly Mutated Genes (SMG). |
| GET | /Analyses/CopyNumber/Genes/All | |
| GET | /Analyses/CopyNumber/Genes/Focal | |
| GET | /Analyses/CopyNumber/Genes/Thresholded | |
| GET | /Analyses/CopyNumber/Genes/Amplified | Retrieve Gistic2 significantly amplified genes results. |
| GET | /Analyses/CopyNumber/Genes/Deleted | |
| GET | /Analyses/Reports | |
| GET | /Analyses/Summary | |

**Samples** Fine grained retrieval of sample-level data    Show/Hide | List Operations

| GET | /Samples/mRNASeq |
| GET | /Samples/miRSeq |
| GET | /Samples/ClinicalTier1 |

**Archives** : Bulk retrieval of data or analysis pipeline results, as compressed archives

Show/Hide | List Operations

| GET | /Archives/StandardData |

**Metadata** : Retrieve disease, sample, and datatype descriptions, sample counts, and more

Show/Hide | List Operations | Expand

| GET | /Metadata/Counts | |
| GET | /Metadata/Cohorts | Retrieve map of cohort abbreviation |
| GET | /Metadata/Cohort/{cohort} | Retrie |
| GET | /Metadata/Platforms | Retrieve map of platform code |

# Interactive Docs

*learn APIs and explore data
by playing in real time
instead of cut/paste from static HTML or PDF*

*automatically generated & updated
as API and database evolve*

**GET** /Samples/mRNASeq

## Implementation Notes

This service returns sample-level log2 mRNASeq expression values. Results may be filtered by gene, cohort, barcode, sample type or characterization protocol, but at least one gene OR barcode must be supplied.

## Parameters

| Parameter | Value | Description | Parameter Type | Data Type |
|---|---|---|---|---|
| format | json (default) | Format of result. | query | string |
| gene | egfr | Comma separated list of gene name(s). | query | string |
| cohort | ACC BLCA BRCA CESC | Narrow search to one or more TCGA disease cohorts from the scrollable list. | query | string |
| tcga_participant_barcode | | Comma separated list of TCGA participant barcodes (e.g. TCGA-GF-A4EO). | query | string |
| sample_type | NB NT TAM TAP | Narrow search to one or more TCGA sample types from the scrollable list. | query | string |
| protocol | RPKM RSEM | Narrow search to one or more sample characterization protocols from the scrollable list. | query | string |

*choices clearly enumerated*

**Perform Query**   Hide Response

**Request URL**

```
http://firebrowse.org:8000/api/v1/Samples/mRNASeq?format=json&gene=egfr&page=1&page_size=250&sort_by=gene
```

*Proper RESTful call is ASSEMBLED FOR YOU*

*Results returned in multiple formats*

```
{
    "cohort": "ACC",
    "expression_log2": 7.59666610237019,
    "gene": "EGFR",
    "geneID": 1956,
    "protocol": "RSEM",
    "sample_type": "TP",
    "tcga_participant_barcode": "TCGA-OR-A5J1",
    "z-score": -0.40056053472322
},
{
    "cohort": "ACC",
    "expression_log2": 6.98214823852598,
    "gene": "EGFR",
    "geneID": 1956,
    "protocol": "RSEM",
    "sample_type": "TP",
    "tcga_participant_barcode": "TCGA-OR-A5J2",
    "z-score": -0.572210443678677
},
```

| tcga_participant_barcode | gene | expression_log2 | z-score | cohort | sample_type | |
|---|---|---|---|---|---|---|
| TCGA-OR-A5J1 | EGFR | 7.59666610237 | -0.400560534723 | ACC | TP | RSEM |
| TCGA-OR-A5J2 | EGFR | 6.98214823853 | -0.572210443679 | ACC | TP | RSEM |
| TCGA-OR-A5J3 | EGFR | 9.31231960446 | 0.729969055244 | ACC | TP | RSEM |
| TCGA-OR-A5J5 | EGFR | 8.50495520815 | 0.0333590221281 | ACC | TP | RSEM |
| TCGA-OR-A5J6 | EGFR | 8.5592941021 | 0.0690092698339 | ACC | TP | RSEM |
| TCGA-OR-A5J7 | EGFR | 8.64932911891 | 0.131115969294 | ACC | TP | RSEM |
| TCGA-OR-A5J8 | EGFR | 8.06454015357 | -0.210987070006 | ACC | TP | RSEM |
| TCGA-OR-A5J9 | EGFR | 6.63334692474 | -0.641628460792 | ACC | TP | RSEM |
| TCGA-OR-A5JA | EGFR | 9.05879837786 | 0.468028706825 | ACC | TP | RSEM |
| TCGA-OR-A5JB | EGFR | 8.50794128032 | 0.0352834298625 | ACC | TP | RSEM |
| TCGA-OR-A5JC | EGFR | 7.55685241318 | -0.414030877529 | ACC | TP | RSEM |
| TCGA-OR-A5JD | EGFR | 6.25656347946 | -0.699966368647 | ACC | TP | RSEM |
| TCGA-OR-A5JE | EGFR | 6.16656683008 | -0.711787657396 | ACC | TP | RSEM |
| TCGA-OR-A5JF | EGFR | 8.56235233966 | 0.0710558865356 | ACC | TP | RSEM |
| TCGA-OR-A5JG | EGFR | 8.96827107766 | 0.385101741143 | ACC | TP | RSEM |
| TCGA-OR-A5JI | EGFR | 7.05755857856 | -0.554865718674 | ACC | TP | RSEM |
| TCGA-OR-A5JJ | EGFR | 6.64321260426 | -0.639886855174 | ACC | TP | RSEM |

*JSON for computers/programmers*

*TSV, CSV for scientists, algorithms*

# Even Easier in Python, R, and UNIX

## fbget

- Low-level Python bindings: 1-1 with RESTful api
- Higher-level interface, for easy/common bioinformatics
- UNIX command line interface, too
- Automatically generated, easily synched with RESTful API
- Flexible, copiously documented and tested
- BSD-style open source license

Download

## FireBrowseR : bindings for R

https://github.com/mariodeng/FirebrowseR

# fbget : low level interface

```python
python>   import firebrowse
python>   print  firebrowse.Samples().mRNASeq(gene="egfr", cohort="ucs")
{
   "mRNASeq": [
      {
         "cohort": "UCS",
         "expression_log2": 7.06162500904694,
         "gene": "EGFR",
         "geneID": 1956,
         "protocol": "RSEM",
         "sample_type": "TP",
         "tcga_participant_barcode": "TCGA-QN-A5NN",
         "z-score": -0.598993525060403
      },
      ...
```

4 classes, one per API category:
*Samples, Analyses,
Archives, Metadata*

N methods per class, matching
RESTful API; each defaults
to returning 1 page, in JSON

# fbget : high level interface

```
python>    import fbget
python>    print fbget.mrnaseq("egfr", cohort="ucs")

tcga_participant_barcode            gene    expression_log2 z-score cohort
TCGA-QN-A5NN      EGFR      7.06162500905    -0.59899352506  UCS        TP
TCGA-QM-A5NM      EGFR      8.16734387649    -0.298443593752 UCS        TP
TCGA-NG-A4VW      EGFR      8.93092623547    0.0932667888031 UCS        TP
```

- Simpler, e.g. objects do not need to be instantiated
- Intuitive defaults for common bioinformatic use cases
- Transparently iterates:
  - ✓ To retrieve all pages of results in 1 call
  - ✓ In TSV format

# fbget : UNIX CLI interface

```
linux%    fbget mrnaseq egfr cohort=ucs

tcga_participant_barcode          gene      expression_log2 z-score cohort
TCGA-QN-A5NN      EGFR      7.06162500905      -0.59899352506  UCS        TP
TCGA-QM-A5NM      EGFR      8.16734387649      -0.298443593752 UCS        TP
TCGA-NG-A4VW      EGFR      8.93092623547      0.0932667888031 UCS        TP
```

Because sometimes even writing just a
couple of lines of Python takes too long

# Example: quickly list patients

**All of TCGA**

```
linux%   fbget patients

tcga_participant_barcode    date        cohort
TCGA-PK-A5H9      2015-04-02 00:00:00 ACC
TCGA-PA-A5YG      2015-04-02 00:00:00 ACC
TCGA-OR-A5JD      2015-04-02 00:00:00 ACC
TCGA-P6-A5OF      2015-04-02 00:00:00 ACC
TCGA-P6-A5OG      2015-04-02 00:00:00 ACC
```

**Or just GBM**

```
linux%   fbget patients   cohort=gbm

tcga_participant_barcode    date        cohort
TCGA-19-4065      2015-04-02 00:00:00 GBM
TCGA-81-5911      2015-04-02 00:00:00 GBM
TCGA-81-5910      2015-04-02 00:00:00 GBM
TCGA-12-1089      2015-04-02 00:00:00 GBM
```

**This can be enhanced to yield platform data matrix, like AWG freeze list**

# fbget Documentation

Docs for almost all class methods and functions can also
be obtained by invoking the function with zero arguments.

```
python>  fbget.mrnaseq()

mrnaseq() call has missing/None arg value(s), need at least one of: gene OR barcode
Help on function mrnaseq in module fbget:

mrnaseq(gene=None, barcode=None, **kwargs)

    High level wrapper for the FireBrowse Samples.mRNASeq method.
    By default it returns ALL pages of data, in TSV format. ▪ ▪ ▪
```

Better than an inscrutable stack trace, don't you think?

# Same is true on UNIX command line

```
linux%  fbget mrnaseq

mrnaseq() call has missing/None arg value(s), need at least one of: gene OR barcode
Help on function mrnaseq in module firebrowse.fbget:

mrnaseq(gene=None, barcode=None, **kwargs)
    High level wrapper for the FireBrowse Samples.mRNASeq method.
    By default it returns ALL pages of data, in TSV format.

    This service returns sample-level log2 mRNASeq expression
    values. Results may be filtered by gene, cohort, barcode,
    sample type or characterization protocol, but at least one
    gene OR barcode must be supplied.

    For more details consult the interactive documentation at
        http://firebrowse.org/api-docs/#!/Samples
    OR use help(param_values) to see the range of values accepted
    for each parameter, the defaults for each (if any), and the
    degrees of optionality/requiredness offered by the API.

    Parameters:
        format  (str)  Format of result.
        gene  (str)  Comma separated list of gene name(s).
        cohort  (str)  Narrow search to one or more TCGA disease cohorts.
        barcode  (str)  Comma separated list of TCGA participant barcodes (e.g. TCGA-GF-A4EO).
        sample_type  (str)  Narrow search to one or more TCGA sample types.
        protocol  (str)  Narrow search to one or more sample characterization protocols.
        page  (int)  Which page (slice) of entire results set should be returned.
        page_size  (int)  Number of records per page of results.  Maximum is 2000.
        sort_by  (str)  Which column in the results should be used for sorting paginated results?
```

Docs obtained by invoking functions with zero arguments

# Examples Embedded Directly in Tool

```
linux%  fbget --examples

    # Every line of these examples can be cut and directly pasted to your
    # UNIX-like command line.  Comments will be ignored, while everything
    # not beginning with the # comment character will be executed, as long
    # as fbget is in your $PATH

    # Get the RNASeq expression level of the POLE gene, for all TCGA samples
    # (both tumors and normals, in RSEM form, saved to file)
    fbget --outfile=fbget-test-pole.tsv mrnaseq pole

    # Similar query, but constrained to just the DLBC disease cohort
    fbget mrnaseq pole cohort=dlbc

    # Now constrained to single patient, and showing case insensitivity
    fbget mrnaseq pOlE baRcOdE=TCGA-RQ-A6JB

    # What is the DLBC cohort, anyway?
    fbget cohort dlbc
    # DLBC      Lymphoid Neoplasm Diffuse Large B-cell Lymphoma

    # List all the disease cohorts offered by FireBrowse (note that aggregate
    # cohorts like COADREAD,KIPAN,GBMLGG,STES are not available at the TCGA DCC)
    fbget cohorts

    # Display help (docstring) for the function which retrieves clinical data
    fbget help clinical
```

# Fin