

Exploring the CCLE DB to identify cell lines with high expression of gene of interest

By Ida Shinder

12/11/17

Download RNAseq file

<https://portals.broadinstitute.org/ccle/data>

[Home](#) [About](#) [Data](#) [Contact](#)

Search Genes/Cell Lines

Search

Logout: ida.shinder

Browse Datasets

Current Data

Filename	Date	Description	Download
CCLE_RNAseq_081117.rpkм.gct	15-Aug-2017	CCLE RNAseq gene expression data (RPKM)	Download 577MB
CCLE_RNAseq_081117.reads.gct	15-Aug-2017	CCLE RNAseq gene expression data (read count)	Download 165MB
ccle2maf_081117.txt	15-Aug-2017	Merged mutation calls (coding region, germline filtered)	Download 250MB

```
[shinderii@helix ~]$ cd /data/shinderii/experiments/2017_12_03_CCLE/Download/  
[shinderii@helix Download]$ wget https://data.broadinstitute.org/ccle/CCLE_RNAseq_081117.rpkм.gct
```

Explore file format/structure

```
[shinderii@helix Download]$ head -n 3 CCLE_RNAseq_081117.rpkm.gct
```

```
#1.2
```

```
56318 1019
```

Name	Description	22RV1_PROSTATE	2313287_STOMACH	253JBV_URINARY_TRACT	253J_URINARY_TRACT	42MGBA_CENTRAL_NERVOUS_SYSTEM	5637_URINARY_TRACT	59
M_OVARY	639V_URINARY_TRACT	647V_URINARY_TRACT	697_HAEMATOPOIETIC_AND_LYMPHOID_TISSUE	769P_KIDNEY	7860_KIDNEY	8305C_THYROID	8505C_THYROID	
8MGBA_CENTRAL_NERVOUS_SYSTEM	A101D_SKIN	A1207_CENTRAL_NERVOUS_SYSTEM	A172_CENTRAL_NERVOUS_SYSTEM	A204_SOFT_TISSUE	A2058_SKIN	A253_SALI		
VARY_GLAND	A2780_OVARY	A375_SKIN	A3KAW_HAEMATOPOIETIC_AND_LYMPHOID_TISSUE	A427	A498_KIDNEY	A4FUK_HAEMATOPOIETIC_AND_LYMPHOID_TISSUE		

Header starts on the third line. When importing to pandas, should skip rows 0 and 1:

```
1 def RNAseq_CCLE(file, gene_list, hgene, order_by):
2     ''' INPUTS
3         file: contains the RNAseq file in gct format
4         gene_list: List of genes of interest (ENSG ID) that includes housekeeping genes
5         hgene: housekeeping gene (to obtain gene/hgene ratio)
6         order_by: gene of importance by which to ascend the gene/hgene ratio
7
8     OUTPUTS:
9         DataFrame that contains cell lines as rows and gene/hgene ratio as columns
10    '''
11
12    #Import file into pandas dataframe
13    df = pd.read_csv(cwd+ '/Download/' + file, skiprows=(0,1), sep = '\t', index_col='Name')
14    df['Name'] = df.index.str.split('.').str.get(0) #The ENSG version is different in the different RNAseq files
15    df = df.set_index('Name')
16
```

Filter by list, add ratio of gene/hgene, reshape df

```
1 def RNaseq_CCLE(file, gene_list, hgene, order_by):
2     ''' INPUTS
3         file: contains the RNaseq file in gct format
4         gene_list: List of genes of interest (ENSG ID) that includes housekeeping genes
5         hgene: housekeeping gene (to obtain gene/hgene ratio)
6         order_by: gene of importance by which to ascend the gene/hgene ratio
7
8     OUTPUTS:
9         DataFrame that contains cell lines as rows and gene/hgene ratio as columns
10    '''
11
12    #Import file into pandas dataframe
13    df = pd.read_csv(cwd+ '/Download/' + file, skiprows=(0,1), sep = '\t', index_col='Name')
14    df['Name'] = df.index.str.split('.').str.get(0) #The ENSG version is different in the different RNaseq files
15    df = df.set_index('Name')
16
17    #Create a new dataframe filtered by genes of interest
18    fdataDF = df.loc[gene_list]
19    fdataDF = fdataDF.set index('Description')
20    trans = fdataDF.transpose()
21
22    #Check for NaN entries (just in case...)
23    if trans[trans.isnull().any(axis=1)].shape[0] != 0:
24        print('WARNING: Transposed table has null values')
25
26    #Add ratios of gene/hgene
27    for i in list(trans):
28        trans[i+'/' + hgene + '_ratio'] = trans[i]/trans[hgene]
29    del trans[i]
30    del trans[hgene+'/' + hgene + '_ratio']
31    del trans['GAPDH/ACTB_ratio']
32
33    #Order data frame by gene of interest
34    sort_by = order_by+'/' + hgene + '_ratio'
35    trans = trans.sort_values(by=sort_by, ascending=0)
36
37    return trans
```

Resulting dataframe structure

```
#Run RNAseq_CCLE function and assign to "dataframe" variable  
dataframe = RNAseq_CCLE('CCLE_RNAseq_081117.rpkm.gct', gene_list, 'ACTB', 'IL10RA')
```

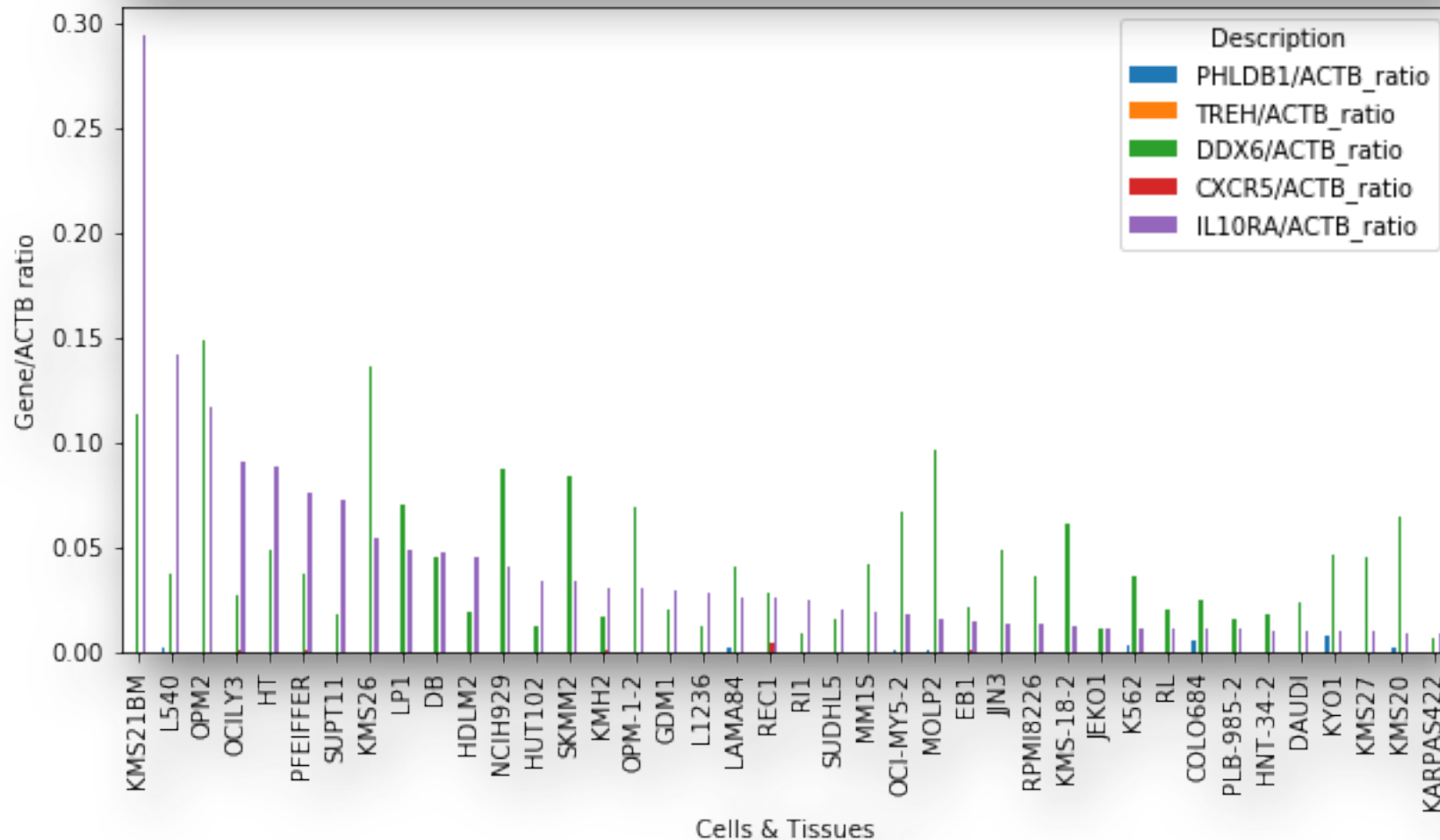
Description	PHLDB1/ACTB_ratio	TREH/ACTB_ratio	DDX6/ACTB_ratio	CXCR5/ACTB_ratio	IL10RA/ACTB_ratio
KMS21BM_HAEMATOPOIETIC_AND_LYMPHOID_TISSUE	0.000071	0.000000	0.113935	0.000000	0.293404
L540_HAEMATOPOIETIC_AND_LYMPHOID_TISSUE	0.002468	0.000000	0.037803	0.000000	0.142410
OPM2_HAEMATOPOIETIC_AND_LYMPHOID_TISSUE	0.000093	0.000022	0.149335	0.000000	0.117272
OCILY3_HAEMATOPOIETIC_AND_LYMPHOID_TISSUE	0.000219	0.000002	0.027593	0.001966	0.091573
HT_HAEMATOPOIETIC_AND_LYMPHOID_TISSUE	0.000080	0.000000	0.049253	0.000454	0.088919

Function to display and save results

```
1 def RNAseq_graph(dataframe,number_of_CL,filename):
2     '''INPUTS
3         dataframe: the dataframe to graph
4         number_of_CL: number of cell lines to show on the x axis
5         filename: desired file name for the graph
6
7         OUTPUT:
8         Graph that displays the cell line on the x axis and the gene/hgene ratio on the y axis
9     '''
10    #Change directory in order to save the png in the Graphs directory
11    os.chdir(cwd + '/Graphs/')
12
13    df = dataframe.head(number_of_CL)
14
15    df['Cell_line'] = df.index.str.split('_').str.get(0)
16    df = df.set_index('Cell_line')
17
18    #Plot and save df
19    plt.rcParams["figure.figsize"] = [8,10]
20
21    ax = df.plot(kind='bar')
22    ax.set_xlabel("Cells & Tissues")
23    ax.set_ylabel("Gene/" + hgene + ' ratio')
24
25    plt.savefig(filename, bbox_inches='tight')
26    plt.show()
```

Graph dataframe, display results

```
4 #Run RNAseq_graph function to graph the output
5 RNAseq_graph(dataframe,40,'IL10RA_ACTB_ratio_CCLE_RNAseq.png')
```



Past & Future Direction

Past:

- NCI-60 Microarrays
- ENCODE RNAseq

Future:

- FANTOM5 CAGE
- CCLE Microarrays
- Incorporate Copy Number from CCLE