

Homework 5

made by Ilya German

Part 1.

1.

Choose the gene: FAT1.

Find the gene on [IntOGen](#).

On screenshot you can see the types of cancer which corresponds with mutations.

Method signals per Cancer Type

Show entries

Search:

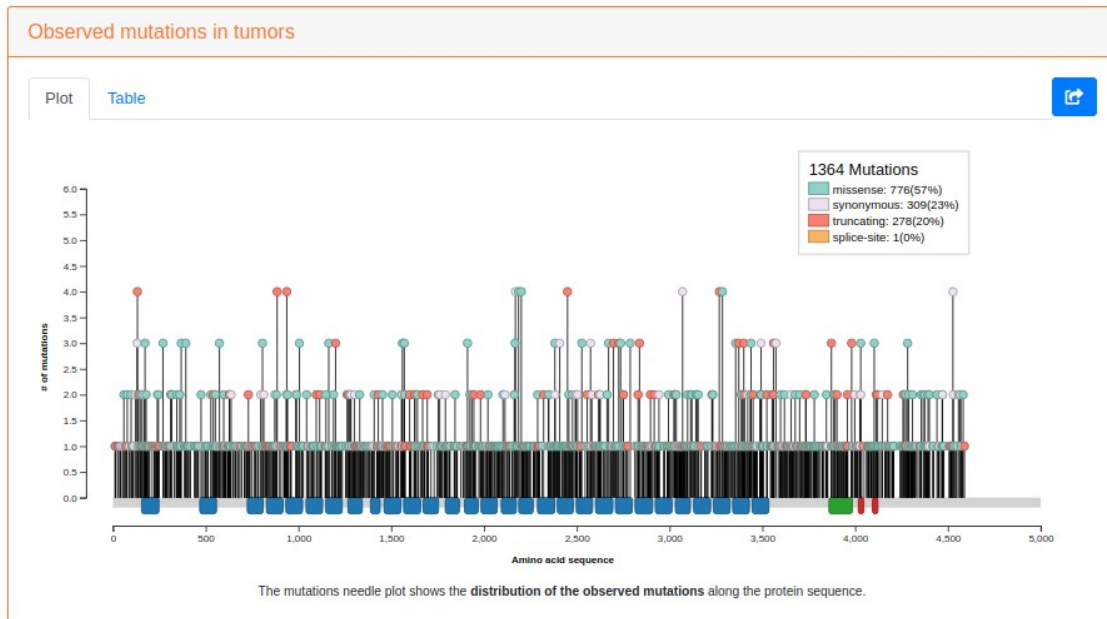
Cancer type	Methods	Samples	Samples (%)
Head and neck squamous cell carcinoma	CBaSE dNdScv MutPanning FML smRegions	131	18.96
Bladder cancer	CBaSE combination dNdScv MutPanning FML	72	8.76
Lung squamous cell carcinoma	CBaSE combination dNdScv MutPanning FML	67	11.47
Endometrial cancer	CBaSE dNdScv	52	8.97
Skin basal cell carcinoma	combination	41	42.71
Cervix squamous cancer	CBaSE dNdScv MutPanning FML	21	6.16
Skin squamous cell carcinoma	CBaSE dNdScv MutPanning	16	32.65
Esophageal cancer	CBaSE dNdScv MutPanning FML	15	1.59
Chronic lymphoblastic leukemia	combination	12	1.65
Ovary cancer	combination	8	1.16
Vulva Cancer	combination	7	31.82
Pancreas adenocarcinoma	combination	5	0.56
Breast adenocarcinoma	smRegions	4	0.15
Prostate adenocarcinoma	combination	3	0.11
Anus cancer	FML	2	11.76
Small intestine cancer neuroendocrine	combination	2	4.44
Hepatic cancer	combination	2	0.12
Medulloblastoma	smRegions	2	0.31
Stomach adenocarcinoma	combination	2	0.28
Uveal melanoma	combination	2	2.5

Showing 1 to 20 of 20 entries

Previous Next

ClustL HotMAPS smRegions Clustered Mutations
CBaSE dNdScv Recurrent Mutations
FML Functional Mutations
MutPanning Tri-nucleotide specific bias
combination Combination

Another screenshot reflects the mutations. There we faces with lots of missense mutations, which refers to point mutations, synonymous mutations, that also refers to point mutations, and truncating mutations that might be both - point or structured mutation. The majority of mutations refers to point-type.



From the next picture we can get, that since 17th page (from 120) the mutations affects only 1 nucleotide, so the majority of truncating mutations are point-type too.

This is only one page with stop_gained consequence, but all of them you can find [here](#).

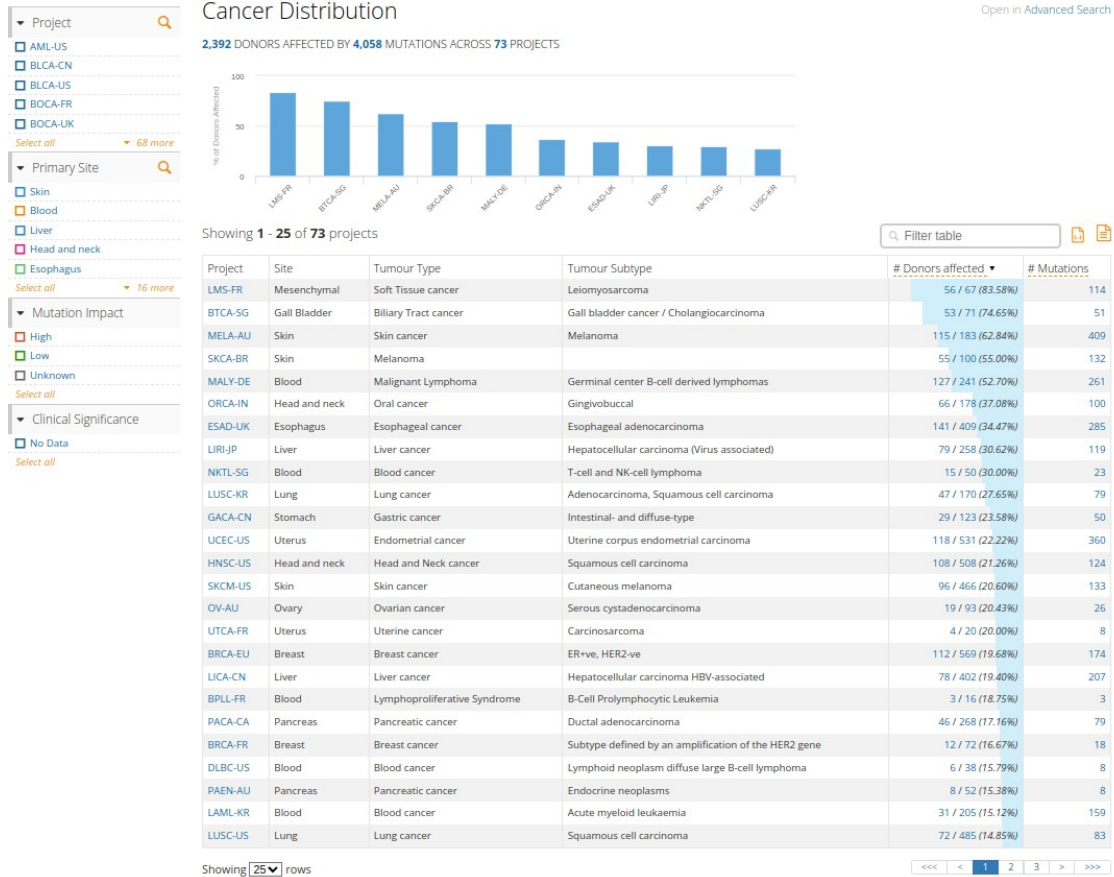
Observed mutations in tumors			
<div> Plot Table </div>		<div> <div>Show 25 entries</div> <div> <div> <div></div> <div>Search:</div> </div> </div> </div>	
Mutation (GRCh38)	Protein Position	Samples	Consequence
4:186600297:G>A	3902	2	stop_gained
4:186603323:G>A	3735	2	stop_gained
4:186603849:G>C	3559	2	stop_gained
4:186603956:G>A	3524	2	stop_gained
4:186606072:G>A	3450	2	stop_gained
4:186606209:G>C	3404	2	stop_gained
4:186614334:G>C	3029	2	stop_gained
4:186617837:G>A	2917	2	stop_gained
4:186617897:G>A	2897	2	stop_gained
4:186618073:G>C	2838	2	stop_gained
4:186618098:G>A	2830	2	stop_gained
4:186618920:G>A	2556	2	stop_gained
4:186619625:G>A	2321	2	stop_gained
4:186620645:G>A	1981	2	stop_gained
4:186620753:T>A	1945	2	stop_gained
4:186621506:G>A	1694	2	stop_gained
4:186621580:G>C	1669	2	stop_gained
4:186621707:G>A	1627	2	stop_gained
4:186628167:G>C	1599	2	stop_gained
4:186633719:C>A	1430	2	stop_gained
4:186636752:C>A	1269	2	stop_gained
4:186636773:G>A	1262	2	stop_gained
4:186639765:G>C	1200	2	stop_gained
4:186663539:G>A	1114	2	stop_gained
4:186663593:G>A	1096	2	stop_gained
Showing 776 to 800 of 1,196 entries		<div> <div>Previous</div> <div>1</div> <div>...</div> <div>31</div> <div>32</div> <div>33</div> <div>...</div> <div>48</div> <div>Next</div> </div>	

2.

Search the ICGC portal for the selected gene [ICGC](#) (need to have VPN).

On screenshot we can see the types of cancer, in which the gene FAT1 mutates. (In the column Tumor Type.)

Again, only one screen is here, all cancer types on [github](#)-page.



Most Frequent Somatic Mutations

[Open in Advanced Search](#)

2. Part 2.

1.

Choose the genome: DO7777.

Download it from [ICGC](https://github.com/ilyagerman52/bioinformatics/raw/main/homework5/simple_somatic_mutation.open.tsv). (I will use my own copy of archive, tha placed on github, because it is easier to code on colab with github instead of ICGC, because there we need to pass info-fields for downloading files.)

```
!wget
```

```
https://github.com/ilyagerman52/bioinformatics/raw/main/homework5/
simple_somatic_mutation.open.tsv
```

```
--2023-03-23 19:50:58--
```

```
https://github.com/ilyagerman52/bioinformatics/raw/main/homework5/
simple_somatic_mutation.open.tsv
```

```
Resolving github.com (github.com)... 140.82.112.3
```

```
Connecting to github.com (github.com)|140.82.112.3|:443... connected.
```

```
HTTP request sent, awaiting response... 302 Found
```

```
Location:
```

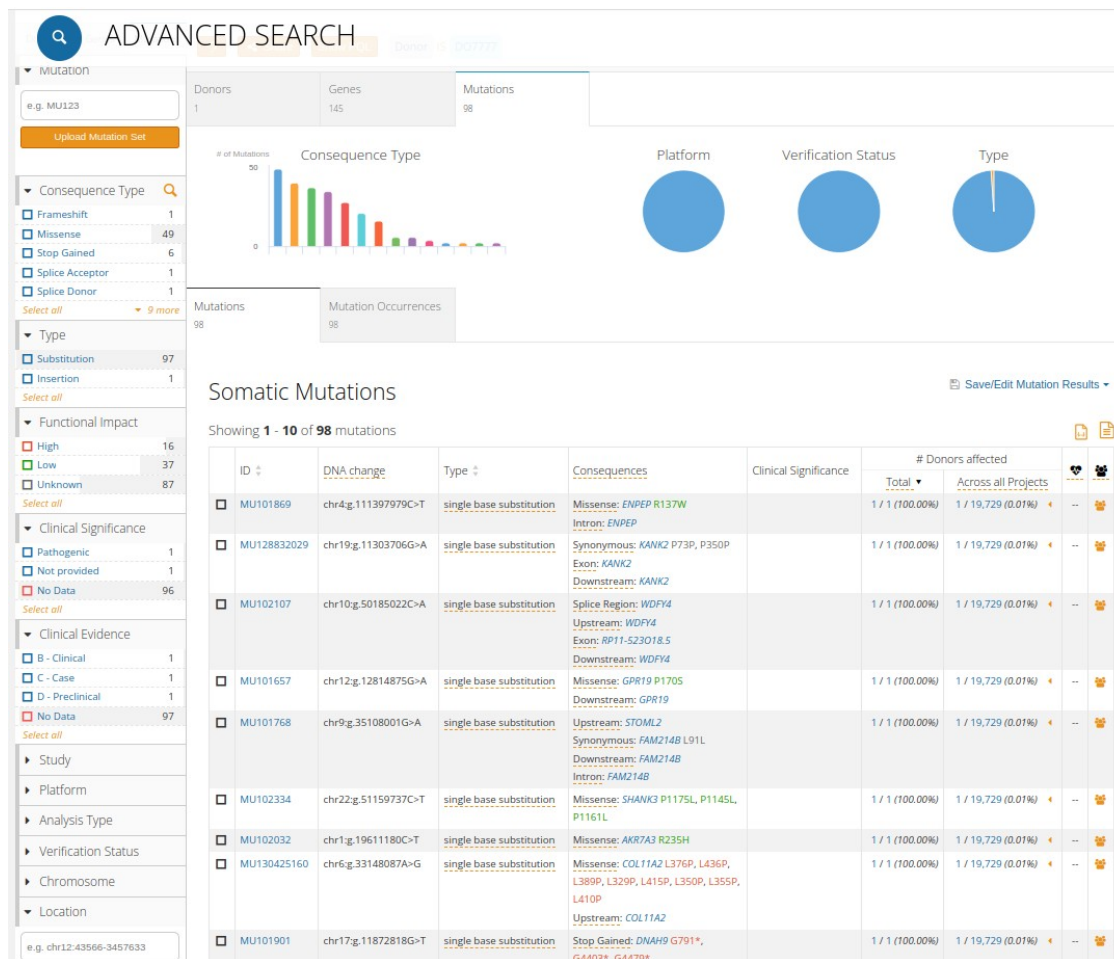
```
https://raw.githubusercontent.com/ilyagerman52/bioinformatics/main/
```

```
homework5/simple_somatic_mutation.open.tsv [following]
--2023-03-23 19:50:58--
https://raw.githubusercontent.com/ilyagerman52/bioinformatics/main/
homework5/simple_somatic_mutation.open.tsv
Resolving raw.githubusercontent.com (raw.githubusercontent.com)...
185.199.108.133, 185.199.109.133, 185.199.110.133, ...
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|
185.199.108.133|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 265136 (259K) [text/plain]
Saving to: 'simple_somatic_mutation.open.tsv'
```

```
simple_somatic_muta 100%[=====>] 258.92K  --.-KB/s    in
0.03s
```


```
2023-03-23 19:50:58 (7.86 MB/s) - 'simple_somatic_mutation.open.tsv'
saved [265136/265136]
```

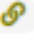
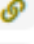
On [page](#) we can find that chosen genome has 98 mutations






For answer the question if any mutations in gene FAT1, we must get the Ensembl ID (for FAT1 it is ENSG00000083857). This is also from [IntOGen](#), the section Gene Details.

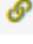

Gene details

FAT1 

Ensembl ID	ENSG00000083857 
Transcript ID	ENST00000441802 
Protein ID	ENSP00000406229

Cancer types where is driver	20
Cohorts where is driver	33
Mutated samples	466
Mutations 	1,396

Mode of action 	Loss of function 
---	--

Known driver	True  
---------------------	--

```
!grep 'ENSG00000083857' simple_somatic_mutation.open.tsv
```

The empty output implies no mutations in FAT1.

Choose another genes. For choosing interesting genes let's see all of ENSG-values of mutated genes. Then find the gene name.

```
!grep -o 'ENSG[0-9]*' simple_somatic_mutation.open.tsv | sort -u
```

```
ENSG00000005238
ENSG00000005302
ENSG00000007174
ENSG00000015479
ENSG00000025293
ENSG00000037280
ENSG00000046889
ENSG00000050730
ENSG00000060656
ENSG00000060718
```

ENSG000000065135
ENSG000000075035
ENSG000000080608
ENSG000000089123
ENSG000000091656
ENSG000000095066
ENSG000000096696
ENSG000000099260
ENSG000000099812
ENSG000000100393
ENSG000000100429
ENSG000000101230
ENSG000000101405
ENSG000000101605
ENSG000000103852
ENSG000000104970
ENSG000000105229
ENSG000000108576
ENSG000000111684
ENSG000000112033
ENSG000000112964
ENSG000000114544
ENSG000000115232
ENSG000000115267
ENSG000000117501
ENSG000000118762
ENSG000000119681
ENSG000000120211
ENSG000000121742
ENSG000000122644
ENSG000000124721
ENSG000000125207
ENSG000000125676
ENSG000000125910
ENSG000000126749
ENSG000000128815
ENSG000000131061
ENSG000000132026
ENSG000000132256
ENSG000000132274
ENSG000000132321
ENSG000000133703
ENSG000000135951
ENSG000000138792
ENSG000000139865
ENSG000000140675
ENSG000000140688
ENSG000000141258
ENSG000000141646
ENSG000000147162

ENSG000000147174
ENSG000000148671
ENSG000000155657
ENSG000000156395
ENSG000000160050
ENSG000000160051
ENSG000000160505
ENSG000000161091
ENSG000000161381
ENSG000000162482
ENSG000000163060
ENSG000000163909
ENSG000000165283
ENSG000000166147
ENSG000000168566
ENSG000000168779
ENSG000000169551
ENSG000000171222
ENSG000000172037
ENSG000000172046
ENSG000000172403
ENSG000000172995
ENSG000000173039
ENSG000000173269
ENSG000000174891
ENSG000000175175
ENSG000000177103
ENSG000000181023
ENSG000000181143
ENSG000000182253
ENSG000000182463
ENSG000000182841
ENSG000000183150
ENSG000000183287
ENSG000000183397
ENSG000000183569
ENSG000000184634
ENSG000000185010
ENSG000000185070
ENSG000000185313
ENSG000000186113
ENSG000000187800
ENSG000000188130
ENSG000000188133
ENSG000000188687
ENSG000000189056
ENSG000000196159
ENSG000000197256
ENSG000000197712
ENSG000000198286

ENSG000000204248
ENSG000000204790
ENSG000000204793
ENSG000000207625
ENSG000000213445
ENSG000000215021
ENSG000000222046
ENSG000000222386
ENSG000000223656
ENSG000000224066
ENSG000000230837
ENSG000000232119
ENSG000000233665
ENSG000000234715
ENSG000000237298
ENSG000000251322
ENSG000000259848
ENSG000000260740
ENSG000000261054
ENSG000000263345
ENSG000000263818
ENSG000000263938
ENSG000000264324
ENSG000000265106
ENSG000000265911
ENSG000000266368
ENSG000000267424
ENSG000000267436
ENSG000000270574
ENSG000000271253
ENSG000000271774
ENSG000000271880
ENSG000000272201
ENSG000000272734
ENSG000000273413

This is Ensembl ID of genes, which are mutated in case of our donor.

I can find some of them, for example let's choose this four:

- ENSG00000075035 - gene [WSCD2](#)
- ENSG000000132321 - gene [IQCA1](#)
- ENSG000000177103 - gene [DSCAML1](#)
- ENSG000000251322 - gene [SHANK3](#)

And two extra genes:

- ENSG000000114544 - gene [SLC41A3](#)
- ENSG000000005238 - gene [ATOSB](#)

Using that grep we find all mutations with ENSG00000075035, then greps by ID (first column with values that start with MU...). Sorting with flag -u efforts to count unique IDs. So, all of these genes have only 1 mutation.

```
print('main 4:\n')
!grep 'ENSG00000075035' simple_somatic_mutation.open.tsv | grep -o
'MU[0-9]*' | sort -u
print('-' * 100)
!grep 'ENSG00000132321' simple_somatic_mutation.open.tsv | grep -o
'MU[0-9]*' | sort -u
print('-' * 100)
!grep 'ENSG00000177103' simple_somatic_mutation.open.tsv | grep -o
'MU[0-9]*' | sort -u
print('-' * 100)
!grep 'ENSG00000251322' simple_somatic_mutation.open.tsv | grep -o
'MU[0-9]*' | sort -u
print('-' * 100)

print('\n\nextras:\n')
!grep 'ENSG00000114544' simple_somatic_mutation.open.tsv | grep -o
'MU[0-9]*' | sort -u
print('-' * 100)
!grep 'ENSG00000005238' simple_somatic_mutation.open.tsv | grep -o
'MU[0-9]*' | sort -u
```

main 4:

MU130337701

MU101543

MU130328708

MU102334

extras:

MU102303

MU101768

The very interesting information is [here](#). For example we can find the mutations MU1330337701 and MU130328708, the corresponds with our mini-research based on greps :)

<input type="checkbox"/>	MU130337701	chr12:g.108634268G>A	single base substitution	Missense: <i>WSCD2</i> R431H		1 / 1 (100.00%)	1 / 19,729 (0.01%)	←	→	🔍
<input type="checkbox"/>	MU130328708	chr11:g.117389195C>T	single base substitution	Missense: <i>DSCAML1</i> R559Q, R289Q		1 / 1 (100.00%)	1 / 19,729 (0.01%)	←	→	🔍

Finally, we can see that the mutation MU130337701 connected with Colon cancer, subtype [Adenocarcinoma](#).

M

MU130337701

GENOTYPE

Summary

Clinical Evidence

Protein

Summary

ID	MU130337701
DNA change	chr12:g.108634268G>A
Type	single base substitution
Reference genome assembly	GRCh37
Allele in the reference assembly	G
Functional Impact	Low

Consequences

Showing 1 Consequences

Open in Advanced Search

Filter table

Gene	AA Change	Consequence	Coding DNA Change	Strand	Transcript(s)
WSCD2	R431H	Missense variant	1292G>A	+	WSCD2-201 WSCD2-002 WSCD2-001 WSCD2-003

Cancer Distribution

MU130337701 AFFECTS 1 DISTINCT DONORS ACROSS 1 CANCER PROJECTS

Showing 1 projects

Open in Advanced Search

Filter table

Project	Site	Tumour Type	Tumour Subtype	# Donors affected
COAD-US	Colorectal	Colon cancer	Adenocarcinoma	1 / 402 (0.25%)