# Tutorial on Mutational Analysis

Mutational processes, DDR mutations and Visualization

Alex Goncearenco

June 16, 2020

Goals:

     1. Identify endogenic and exogenic mutational processes active in the samples using mutational data

     2. Identify mutated DNA damage repair (DDR) genes that may be responsible for the endogenic mutagenesis

Outline:

     1. Mutational data

     2. Software and resources

     3. Running the calculations

     4. Visualizing the results

# Data

**Required (any of the following):**
- GDC / TCGA somatic mutations in MAF format
- MAF or VCF files from MSKCC, ICGC or any other repository
- Or Custom-made MAF (e.g. sequenced cell lines)

**Required:**
- List of DDR genes

**Optional:**
- Annotation of signatures
- Additional sample annotations (e.g. CIMP status)

# Mutation data specifics

- Germline or somatic
- "Simple" mutations vs structural changes and copy number variation
- NGS mutation callers (MuTect 1,2; Strelka, GATK, Somatic Sniper, VarScan2, VarDict, …)
- VCF files – Variant call format
- MAF files – Mutation annotation format (vcf2maf)
- Genome assembly (hg19, hg38)
- One file per sample vs multiple samples in one file

# Software and resources

- MutaGene – analysis of mutational signatures and motifs
  - https://www.ncbi.nlm.nih.gov/research/mutagene/
  - https://pypi.org/project/mutagene/

- COSMIC mutational signatures
  - https://cancer.sanger.ac.uk/cosmic/signatures

- DNA Damage Repair (DDR) genes
  - https://www.mdanderson.org/documents/Labs/Wood-Laboratory/human-dna-repair-genes.html
  - https://gist.github.com/neksa/df619d78cbc42c9a1abab8aa4807bfe7

- Clustergrammer
  - http://amp.pharm.mssm.edu/clustergrammer/

- Python, Jupyter notebook and Anaconda
  - https://www.anaconda.com/products/individual
  - https://github.com/MaayanLab/clustergrammer-widget

# Download mutational data (TCGA)

https://portal.gdc.cancer.gov/repository



cases.project.project_id in ["TCGA-COAD"] and files.access in ["open"] and files.data_format in ["maf"]

| | Access | File Name | Cases | Project | Data Category | Data Format | File Size | Annotations |
|---|---|---|---|---|---|---|---|---|
| 🛒 | 🔓 open | TCGA.COAD.somaticsniper.70835251-ddd5-4c0d-968e-1791bf6379f6.DR-10.0.somatic.maf.gz | 433 | TCGA-COAD | Simple Nucleotide Variation | MAF | 41.62 MB | 85 |
| 🛒 | 🔓 open | TCGA.COAD.varscan.8177ce4f-02d8-4d75-a0d6-1c5450ee08b0.DR-10.0.somatic.maf.gz | 433 | TCGA-COAD | Simple Nucleotide Variation | MAF | 59.93 MB | 85 |
| 🛒 | 🔓 open | TCGA.COAD.muse.70cb1255-ec99-4c08-b482-415f8375be3f.DR-10.0.somatic.maf.gz | 432 | TCGA-COAD | Simple Nucleotide Variation | MAF | 51.14 MB | 85 |
| 🛒 | 🔓 open | TCGA.COAD.mutect.03652df4-6090-4f5a-a2ff-ee28a37f9301.DR-10.0.somatic.maf.gz | 433 | TCGA-COAD | Simple Nucleotide Variation | MAF | 65.59 MB | 85 |

# Download and unpack files



```
gunzip TCGA.COAD.mutect.03652df4-6090-4f5a-a2ff-ee28a37f9301.DR-10.0.somatic.maf.gz
```

# Download metadata (TCGA)



**Download JSON and save as TCGA_COAD_cases.json**

# What is MutaGene?

# Identify mutational processes in mutagene

# Identify mutational processes in mutagene

# What are COSMIC signatures?

https://cancer.sanger.ac.uk/cosmic/signatures



## Proposed aetiology

An endogenous mutational process initiated by spontaneous or enzymatic deamination of 5-methylcytosine to thymine which generates G:T mismatches in double stranded DNA. Failure to detect and remove these mismatches prior to DNA replication results in fixation of the T substitution for C.

## Associated mutation classes and signatures

The activity of SBS1 is closely correlated with the activity of SBS5 within many types of cancer. However, between cancer types, mutation burdens of SBS1 and SBS5 do not clearly correlate consistent with them being due to different underlying processes.

# What are COSMIC signatures?

- '1': 5mC>T
- '2': AID/APOBEC
- '13': AID/APOBEC
- '3': HR DDR defect
- '4': Tobacco
- '10a': Polymerase epsilon
- '10b': Polymerase epsilon
- '14': Polymerase epsilon + MMR
- '15': MMR
- '7b': UV light
- '18': reactive oxygen species
- '24': guanine repair by NER

- '5': Unknown
- '6': DDR MSI
- '19': Unknown
- '20': MMR, POLD1
- '21': MMR
- '26': MMR
- '32': Azathioprine
- '44': MMR
- '46': Sequencing artifact
- '54': Sequencing artifact
- '84': AID

**Download COSMICv3_labels.yml:**
https://gist.github.com/neksa/089af7bea566f76ceabd7337b497bb96

# What are DDR genes?

- DNA Damage Repair (DDR) genes
  - https://www.mdanderson.org/documents/Labs/Wood-Laboratory/human-dna-repair-genes.html

  **Download and save DDR_genes.yml file:**
  https://gist.github.com/neksa/df619d78cbc42c9a1abab8aa4807bfe7

MMR:
  - MSH2
  - MSH3
  - MSH6
  - MLH1
  - PMS2
  - MSH4
  - MSH5
  - MLH3
  - PMS1
  - PMS2P3
  - PMS2L3

PARP:
  - PARP1
  - ADPRT
  - PARP2
  - ADPRTL2
  - PARP3
  - ADPRTL3

Direct-reversal:
  - MGMT
  - ALKBH2
  - ABH2
  - ALKBH3
  - DEPC1

Topo-crosslinks:
  - TDP1
  - TDP2
  - TTRAP

**Consider Missense or Nonsense mutations**

**Do not differentiate between mutations**

# What is Clustergrammer?



- Clustergrammer
  - http://amp.pharm.mssm.edu/clustergrammer/

# What is Anaconda and Jupyter notebook?

- Python, Jupyter notebook and Anaconda
  - https://www.anaconda.com/products/individual
  - https://github.com/MaayanLab/clustergrammer-widget

- Jupyter
    https://jupyter-notebook.readthedocs.io/en/stable/ui_components.html

# What is Anaconda and Jupyter notebook?

# What is Anaconda and Jupyter notebook?

# Setup conda

```
# open terminal
# if your Anaconda is installed you should have "(base)" prefix in your terminal prompt

conda create -n mutations python=3.7 jupyter
conda activate mutations

pip install mutagene
pip install pyyaml
pip install pandas==0.25.1

pip install clustergrammer-widget
jupyter nbextension enable --py --sys-prefix widgetsnbextension
jupyter nbextension enable --py --sys-prefix clustergrammer_widget

mutagene --version

# change directory to where you will be analyzing data and storing files
jupyter notebook
```
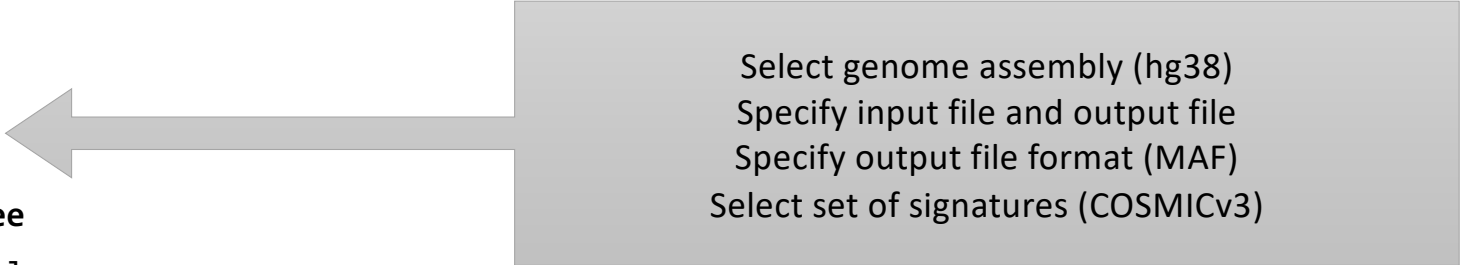
# Setup mutagene package

```
# change directory to where you will be analyzing data and storing files

# choose appropriate genome assembly

# -g hg38 for GRCh38

# -g hg19 for GRCh37


mutagene fetch genome -g hg38
```

# Run MutaGene to identify mutational signatures

```
mutagene identify
        -g hg38
        -i TCGA.COAD.mutect.03652df4-6090-4f5a-a2ff-ee28a37f9301.DR-10.0.somatic.maf
        -o COAD.txt
        -f MAF
        -s COSMICv3
```



Select genome assembly (hg38)
Specify input file and output file
Specify output file format (MAF)
Select set of signatures (COSMICv3)

**# for other options see**

mutagene identify --help

**# mutational decomposition into signatures results are now in COAD.txt file**

(mutations) HG-02113788-LM4:MutationalAnalysisDemo gonceare$ head COAD.txt

| sample | signature | exposure | mutations |
|---|---|---|---|
| TCGA-AA-3966-01A-01D-1981-10 | COSMICv3-SBS15 | 0.223284 | 201 |
| TCGA-AA-3966-01A-01D-1981-10 | COSMICv3-SBS6 | 0.125008 | 113 |
| TCGA-AA-3966-01A-01D-1981-10 | COSMICv3-SBS20 | 0.114956 | 104 |
| TCGA-AA-3966-01A-01D-1981-10 | COSMICv3-SBS54 | 0.0984864 | 89 |

…

# Run jupyter notebook

```
# to run notebook server and open a browser window:

jupyter notebook
```

https://github.com/elnitskilab/MutationalAnalysisDemo