

Signatures of mutational processes in human cancer

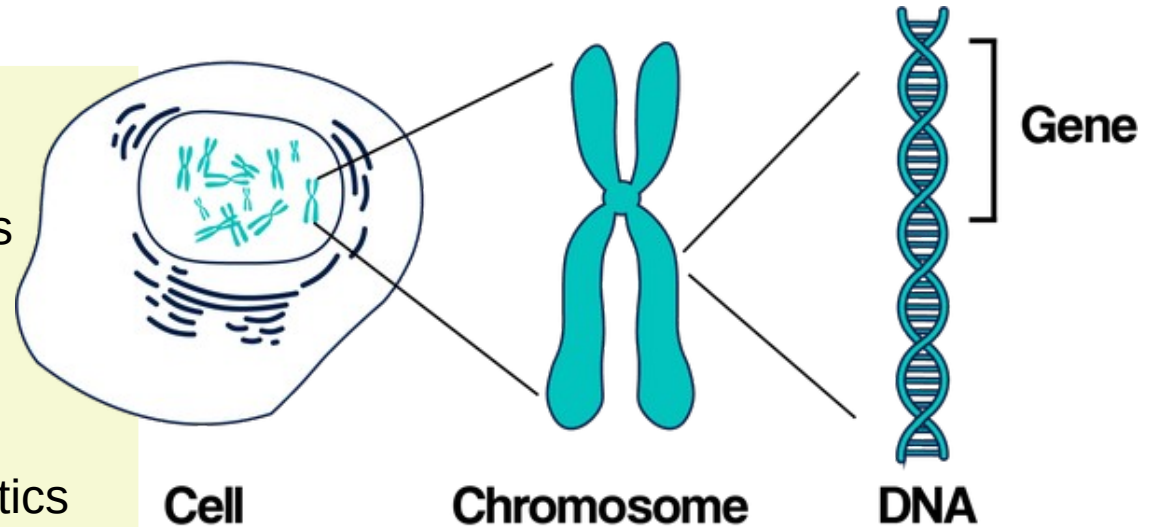
Noel László Plaszkó

Data Science Laboratory
Supervisor: Orsolya Pipek

Introduction: chromosomes, genes

Chromosomes

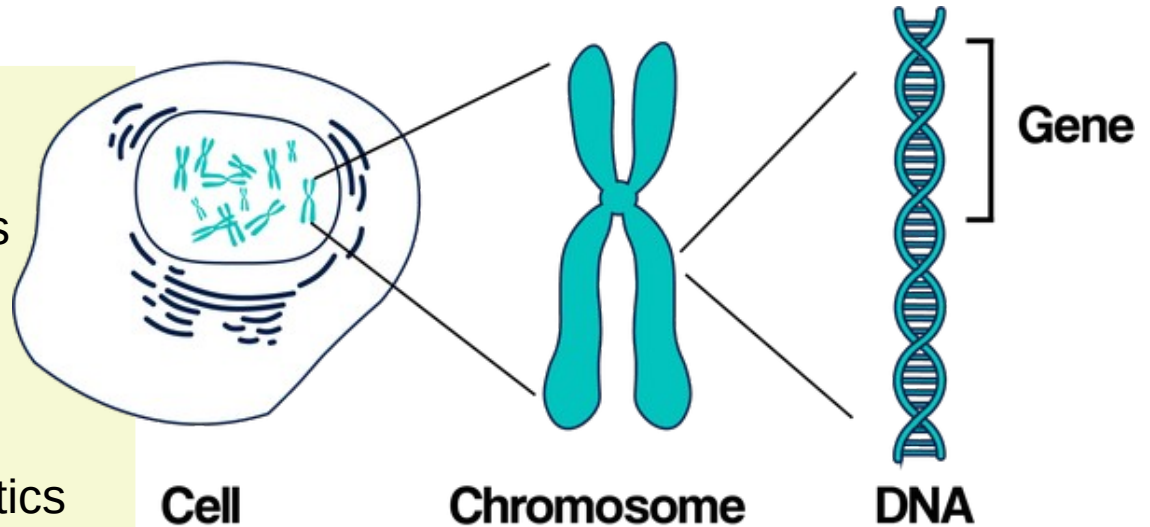
- thread-like structures in cells
 - 46 chromosomes (arranged in 2 sets of 23)
 - one set from your mother, one from your father
- determine physical characteristics
 - contain genes



Introduction: chromosomes, genes

Chromosomes

- thread-like structures in cells
 - 46 chromosomes (arranged in 2 sets of 23)
 - one set from your mother, one from your father
- determine physical characteristics
 - contain genes



Genes

- the basic physical unit of inheritance
- control how your cells work by making proteins
- must have the correct instructions for making its protein
- allows the protein to perform the correct function for the cell

Sorry, something went wrong...

- Genes in cells can mutate
- Due to mutation:
 - Abnormal protein is created
 - Or protein's formation is prevented



Abnormal protein provides different information



This can cause cells to multiply uncontrollably and become cancerous

**Cancer is caused
by somatic mutations of the DNA in the tumor cells.**

DNA sequencing and mutations

DNA sequence:

- the order of the four bases in DNA (adenine, guanine, cytosine, and thymine)

gctagctataggaccgactatcgaaatgctttcctggcgacgcgcgctctggacgtct
atatcgtctgttaaccaagtcggagccgctcagtcataagtaaaacgctacgtttgcttc
tggtcaatgaccagtcgaatttcttatgacgcaatacccttcgtccctgcgaggaactg
gtccagctacaggtacgcggactagtgtatcacgcgacgtgcagcccccagcgcgtatcgg
ccctcgtcgtgcttctgaatcgggttttggatgctgggcgacgtgcagcccccagcgcgtatcgg
ccggcgcggaatccaatgctctagcaagaggggtgccggtgcagcccccagcgcgtatcgg
gcggaatgaacaaaacagtcgaacggcagcgcgttcgtgaggtgcagcccccagcgcgtatcgg
tggcaagtcgtaatatgctctaaaataagtgctctacatattactatgggcccgcgacgc
ggggaagggttgagaaacgaactcggaggggtaggtacgtgaatggcttgaggtgaggt
gttcaagaggtgtaataacgtaacacgttcgaaggggtaggtacgtgaatggcttgaggtgaggt
gcgtcataaacacccatcgctaataacgttcgacctgacacggaggggttttgaggtatgcgg
gccaattaacaccacggggaataattaacgaacaaatccgtgttatctcaagtaggaacc
aacgtctcggctcgtacacgctatcgagagactcgaatgtaaaaccgatacacgaaaagg
actaatagctctactgctcgttgggggtaggtcgaacaaactaggaggtcgaactaagcctg
gcaacgagcgttgctctcgcgactcgggtcgaacaaactaggaggtcgaactaagcctg
actcccgccggaattcgactccattccattaggaacatgaagagcgcgactaaattta

DNA mutation:



Where my story begins...

Somatic mutations for different cancer types

File name	Type of cancer
KIRC.maf	kidney renal clear cell carcinoma
LUAD.maf	lung adenocarcinoma
LUSC.maf	lung squamous cell carcinoma
OV.maf	ovarian cancer
PRAD.maf	prostate adenocarcinoma

Somatic mutations for different cancer types

File name	Type of cancer
KIRC.maf	kidney renal clear cell carcinoma
LUAD.maf	lung adenocarcinoma
LUSC.maf	lung squamous cell carcinoma
OV.maf	ovarian cancer
PRAD.maf	prostate adenocarcinoma

Files have many columns, but only a few will be necessary



For what?!



Different mutational processes (“signatures”) exist



Generate different combinations of mutation types

Somatic mutations for different cancer types

File name	Type of cancer	Number of Samples	Number of mutations
KIRC.maf	kidney renal clear cell carcinoma	235	26245
LUAD.maf	lung adenocarcinoma	561	232492
LUSC.maf	lung squamous cell carcinoma	497	173223
OV.maf	ovarian cancer	142	6174
PRAD.maf	prostate adenocarcinoma	499	36805

Files have many columns, but only a few will be necessary



For what?!



Different mutational processes (“signatures”) exist



Generate different combinations of mutation types

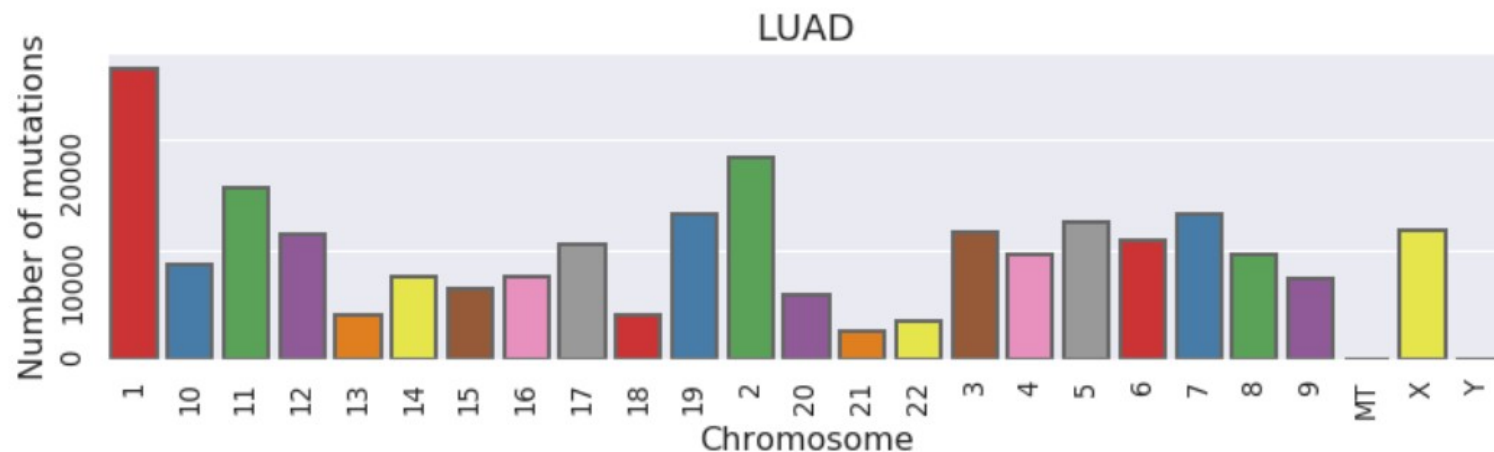
What is in the .maf files?

- Hugo_Symbol
- Entrez_Gene_Id
- Center
- Ncbi_Build
- **Chrom**
- **Start_Position**
- **End_Position**
- Strand
- Variant_Classification
- Variant_Type
- **Reference_Allele**
- Tumor_Seq_Allele1
- **Tumor_Seq_Allele2**
- Dbsnp_Rs
- Dbsnp_Val_Status
- Tumor_Sample_Barcode
- Matched_Norm_Sample_Barcode
- Match_Norm_Seq_Allele1
- Match_Norm_Seq_Allele2
- Tumor_Validation_Allele1
- Tumor_Validation_Allele2
- Match_Norm_Validation_Allele1
- Match_Norm_Validation_Allele2
- Verification_Status
- Validation_Status
- Mutation_Status
- Sequencing_Phase
- Sequence_Source
- Validation_Method
- Score
- Bam_File
- Sequencer
- **Tumor_Sample_UUID**
- Matched_Norm_Sample_UUID
- File_Name
- Archive_Name
- Line_Number



Chrom
Start_Position
End_Position
Reference_Allele
Tumor_Seq_Allele2
Tumor_Sample_UUID

What else do we need?

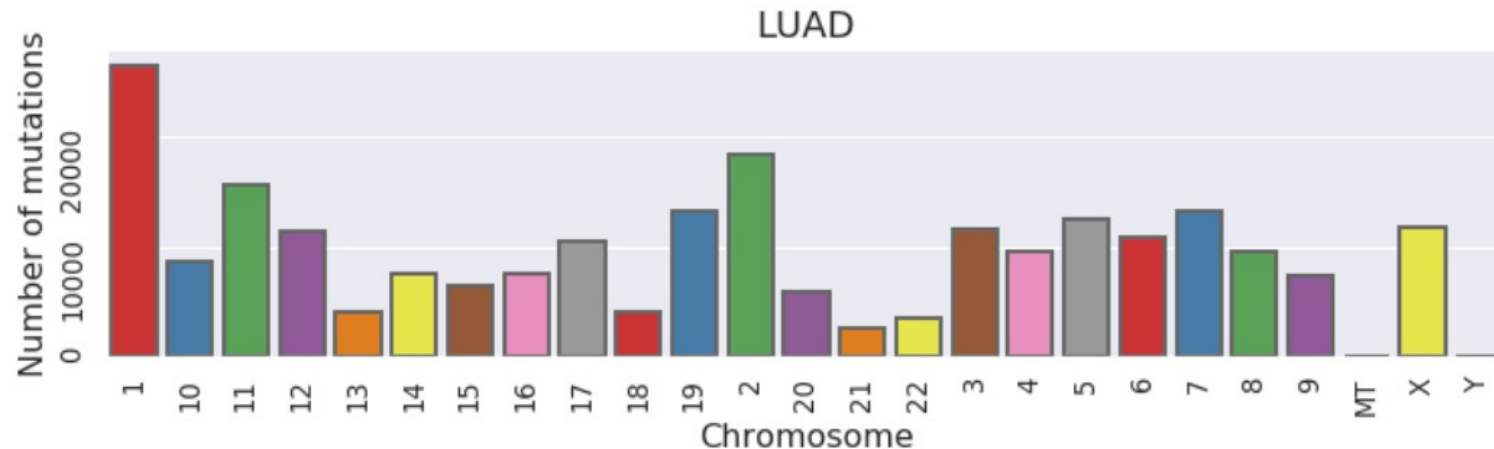


Info from substitution is not enough



Take the neighbours into account

What else do we need?



Info from substitution is not enough



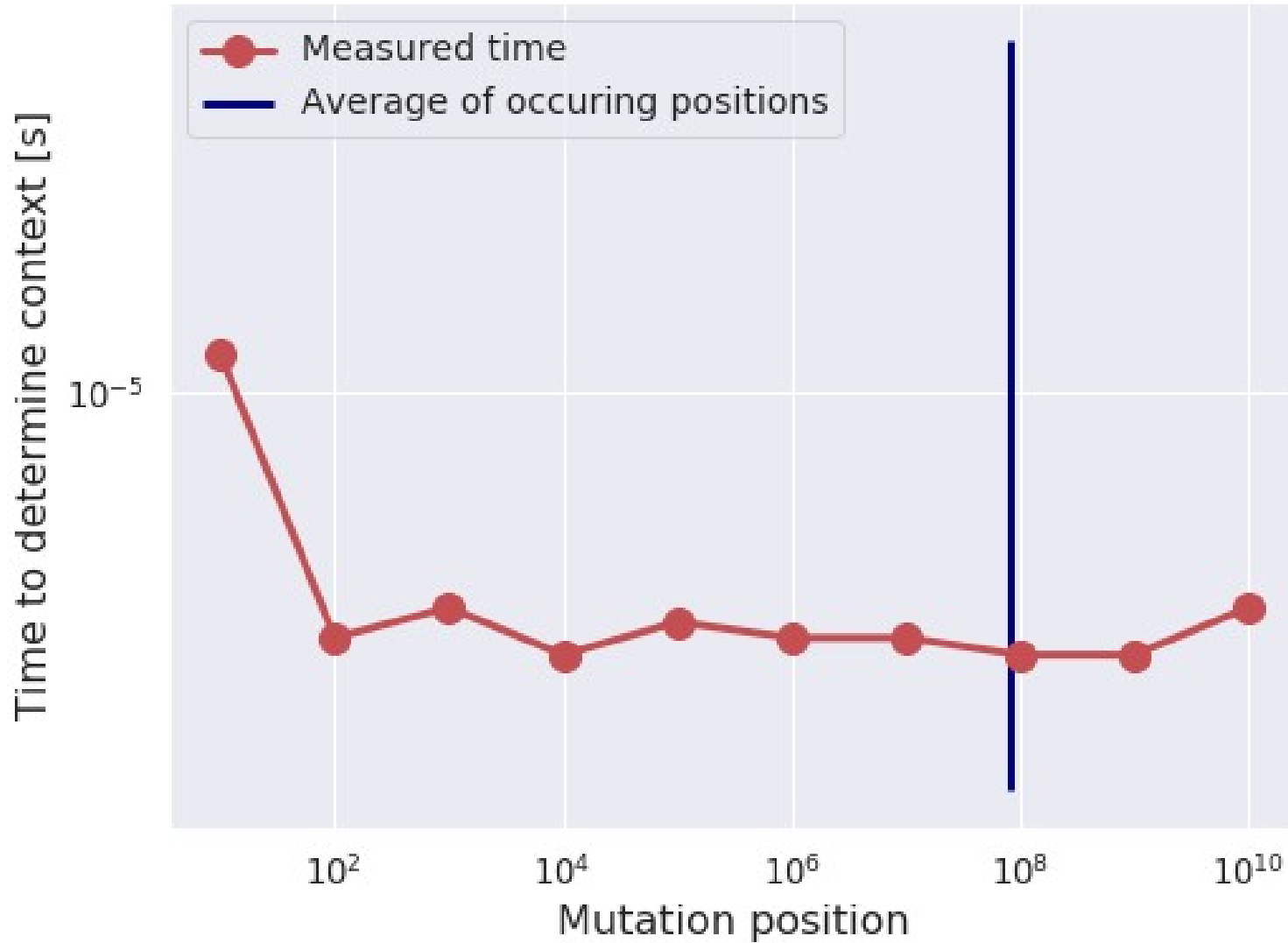
Take the neighbours into account



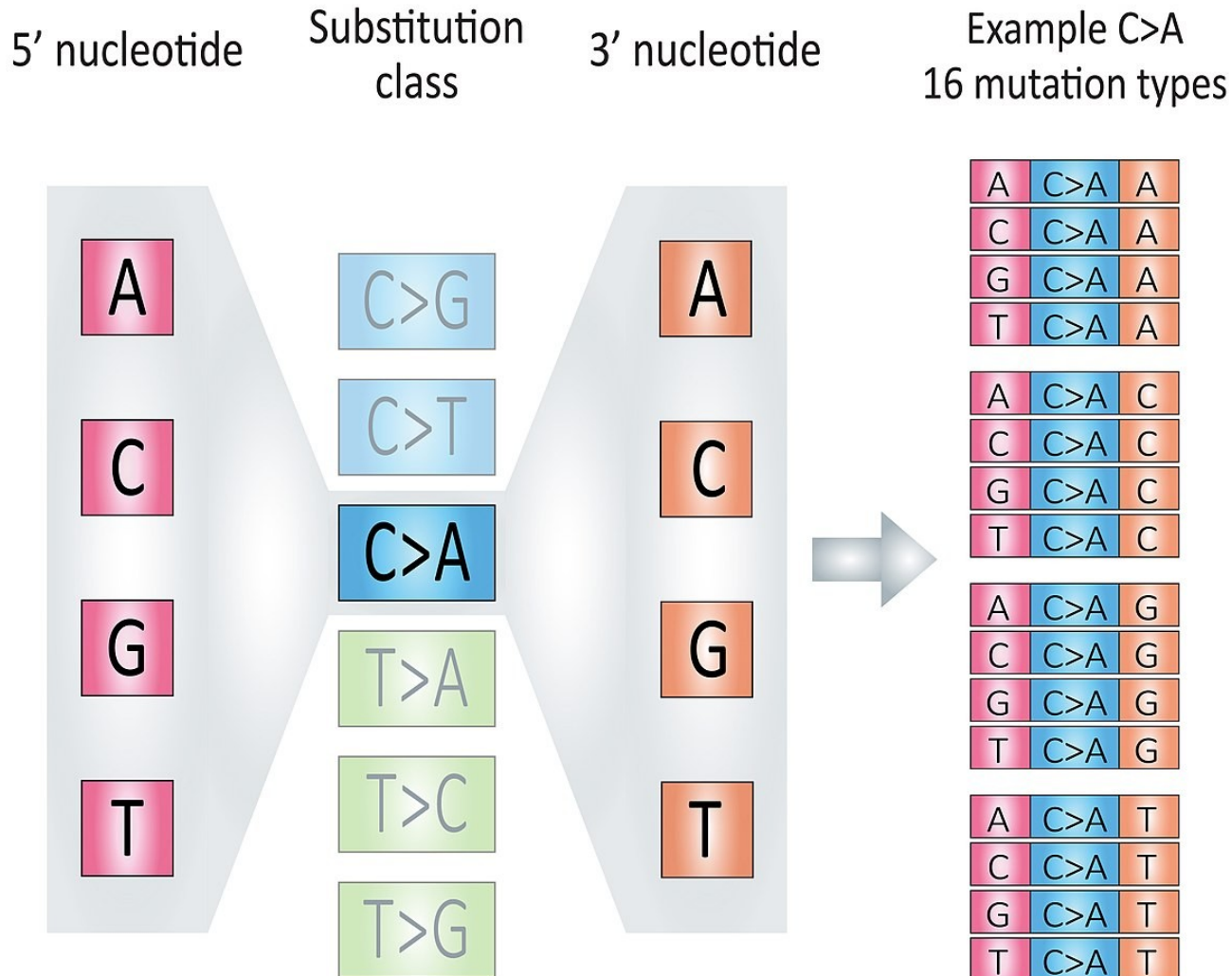
For this purpose **hg19 human reference genome** is used

```
...TGTGTGTGTCCACACTTCCTCTCATGAGAACAG  
CAGGTTGCTTTAGGGCCCACCCTGACAGCCTCGTTC  
TAATACTATGAGGCCAAATACTCACGTTCT...
```

What else do we need?



Mutational catalogs



**Go to the notebook,
to see catalogs...**

Non-negative matrix factorization

Goal: Find two non-negative matrices (W , H) whose product approximates the non-negative matrix X

Objective:

$$\begin{aligned} & 0.5 * ||X - WH||_{loss}^2 \\ & + \alpha_W * l1_{ratio} * n_{features} * ||vec(W)||_1 \\ & + \alpha_H * l1_{ratio} * n_{samples} * ||vec(H)||_1 \\ & + 0.5 * \alpha_W * (1 - l1_{ratio}) * n_{features} * ||W||_{Fro}^2 \\ & + 0.5 * \alpha_H * (1 - l1_{ratio}) * n_{samples} * ||H||_{Fro}^2 \end{aligned}$$

Non-negative matrix factorization

Goal: Find two non-negative matrices (W, H) whose product approximates the non-negative matrix X

Objective:

$$\begin{aligned} & 0.5 * ||X - WH||_{loss}^2 \\ & + \alpha_W * l1_{ratio} * n_{features} * ||vec(W)||_1 \\ & + \alpha_H * l1_{ratio} * n_{samples} * ||vec(H)||_1 \\ & + 0.5 * \alpha_W * (1 - l1_{ratio}) * n_{features} * ||W||_{Fro}^2 \\ & + 0.5 * \alpha_H * (1 - l1_{ratio}) * n_{samples} * ||H||_{Fro}^2 \end{aligned}$$

Hyperparameter tuning:

$\alpha_W = \alpha_H = \alpha$

l1 ratio

number of Components

normalization of X

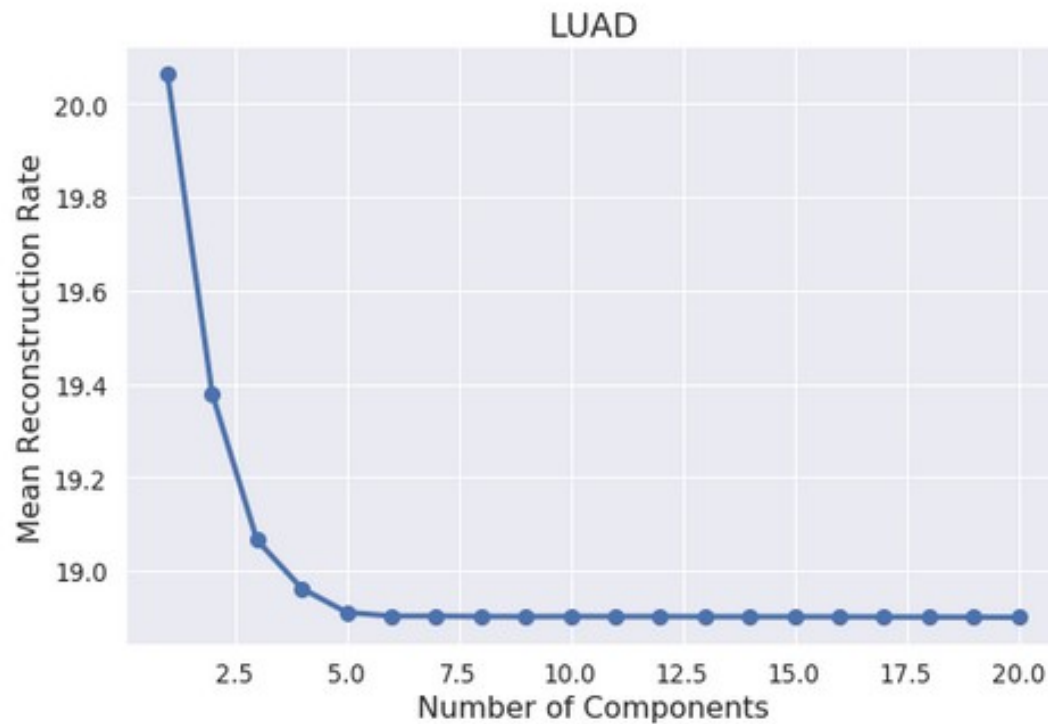
Mutation Max Norm Sample Max Norm Mutation Sum Norm Sample Sum Norm

Non-negative matrix factorization

Choosing the best parameters:

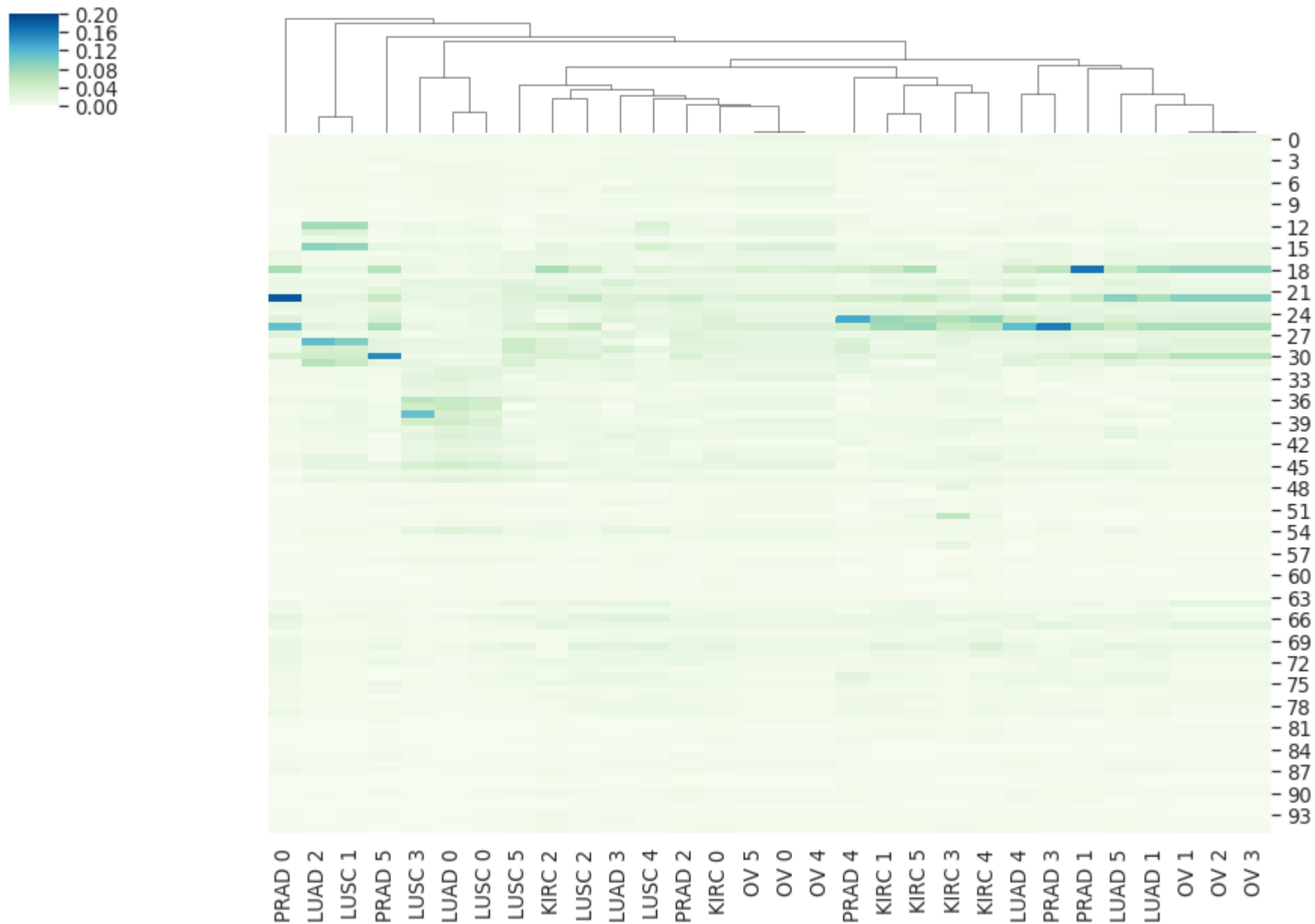
Reconstruction error is calculated

(Frobenius norm of the matrix difference between the training data X and the reconstructed data WH from the fitted model)

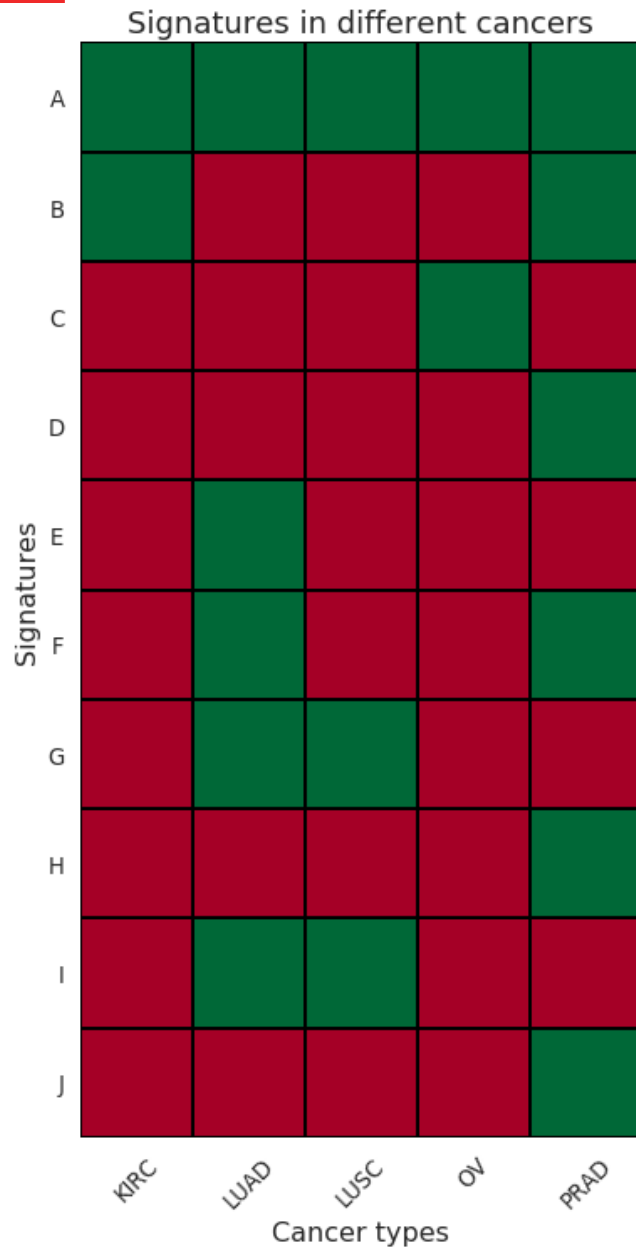


6 number of components are kept → totally 30 signatures

Clustering signatures



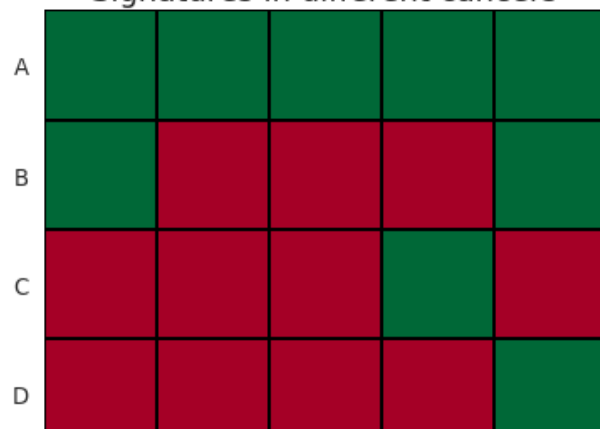
Clustering signatures



- Hierarchical clustering and K-means give similar result

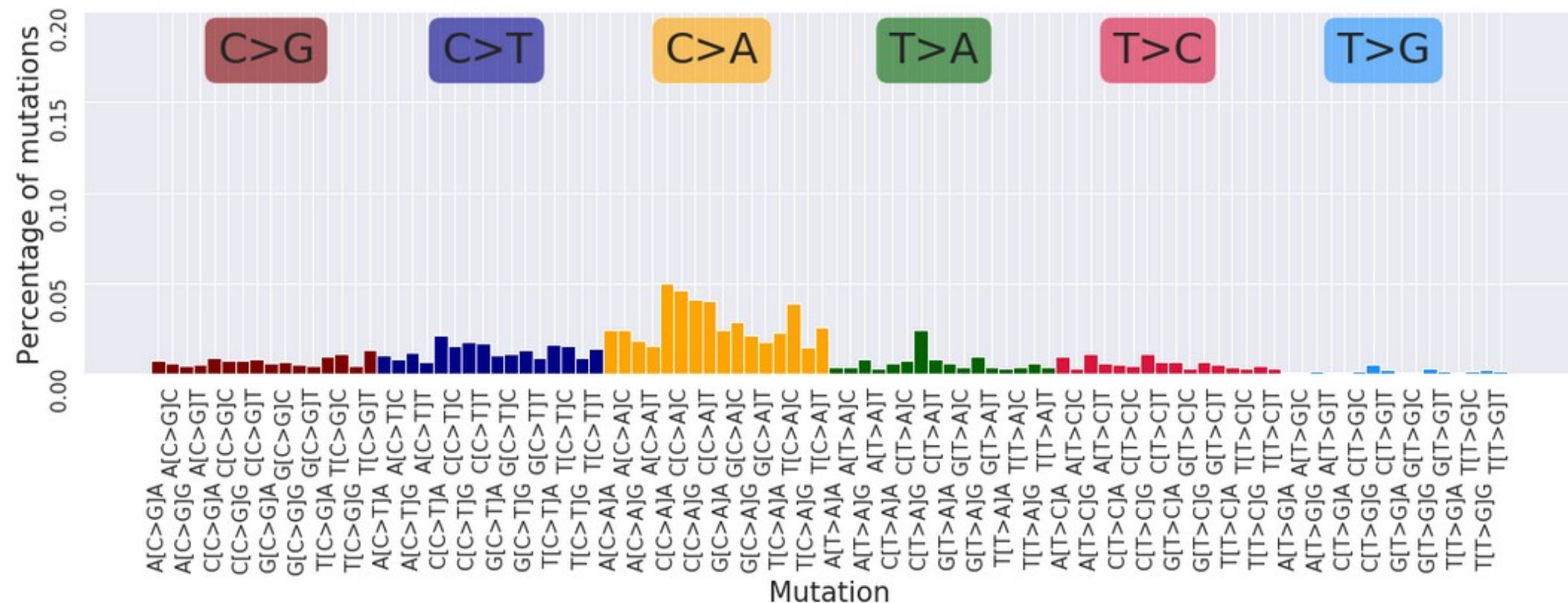
Clustering signatures

Signatures in different cancers



- Hierarchical clustering and K-means give similar result
- **See more figures in the notebook!**

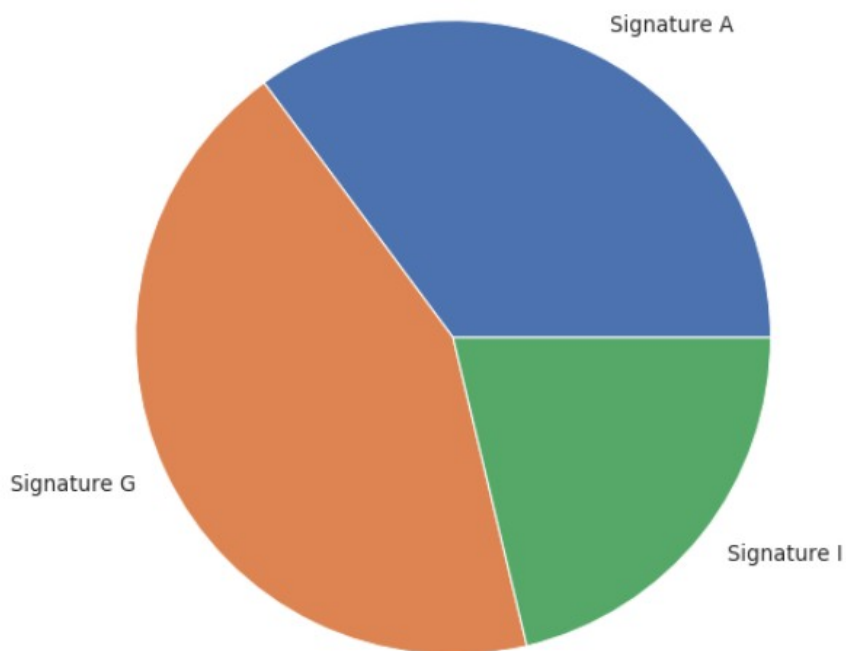
Signature G



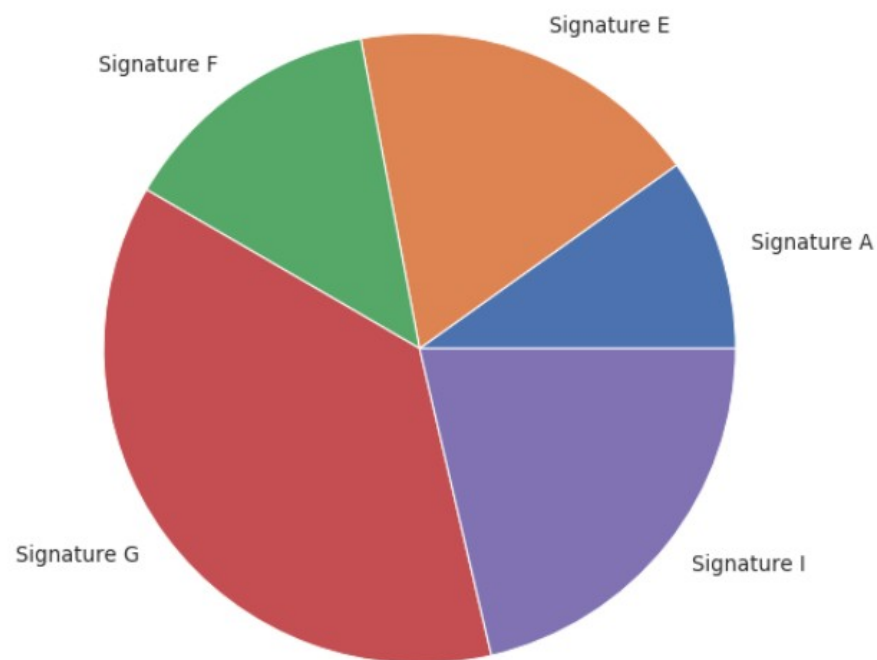
Conclusion, further questions

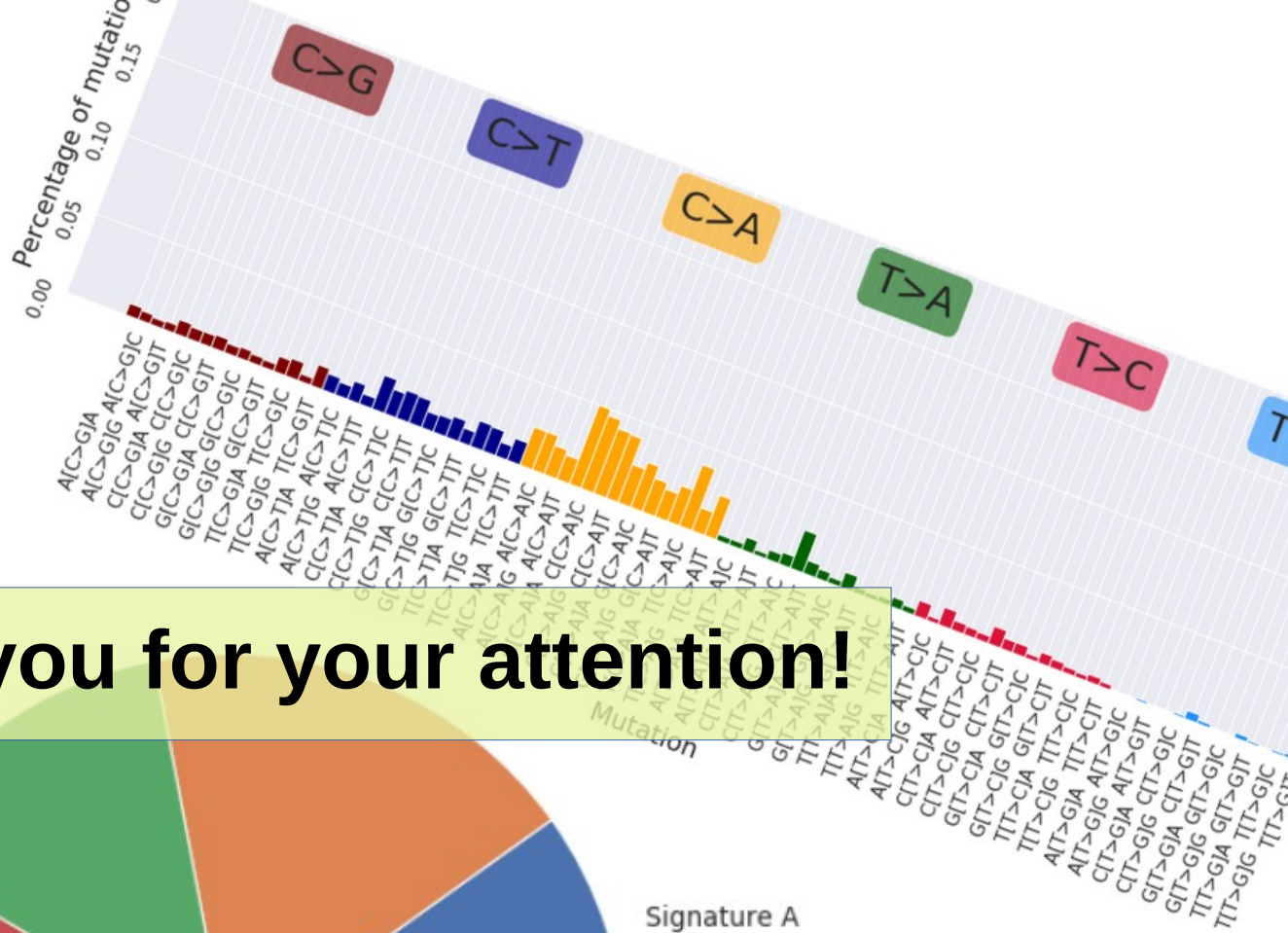
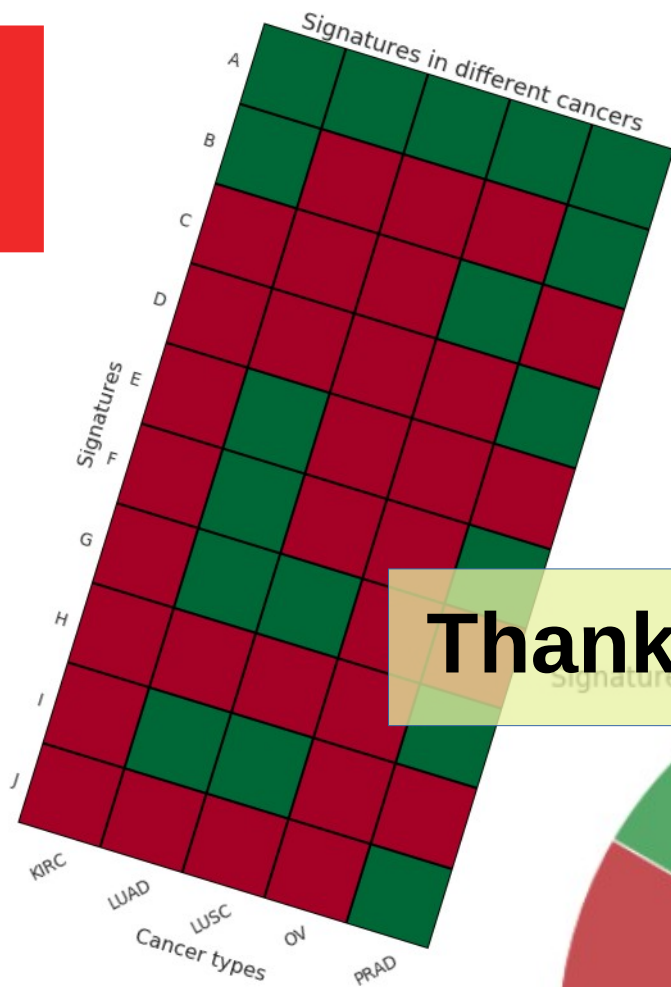
- Mutational signatures are identified
- How are they related to biological processes?
- How are they presented on different samples?

LUSC



LUAD





Thank you for your attention!

