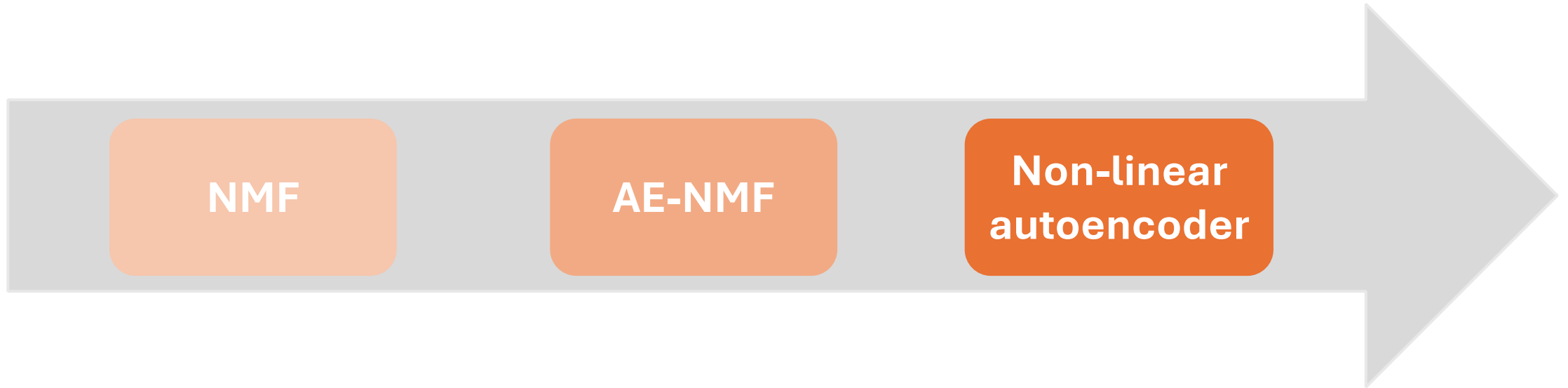# Autoencoders for mutational signature extraction
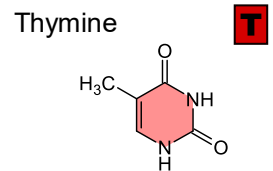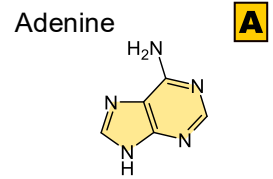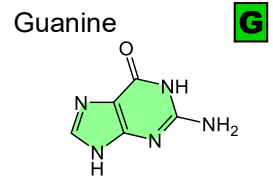
Cortinovis Nicola, Lucas Marta, Paladino Annalisa

UNIVERSITÀ DEGLI STUDI DI TRIESTE

# Introduction

**Goal:** Study mutational signature extraction using autoencoders

# Mutational signatures



Cytosine — C

Guanine — G

Adenine — A

Thymine — T

Nucleobases of DNA

Base pair

**DNA**
Deoxyribonucleic acid

## Single Base Substitutions (SBS)

C > A   C > G   C > T   T > A   T > C   T > G

## Context

AC > AA
AC > AC
AC > AG
AC > AT
CC > AA
CC > AC

96 combinations

# Types of mutational processes

Endogenous

Exogenous



*SBS4 mutational signature linked to tobacco smoking, COSMIC (https://cancer.sanger.ac.uk/signatures/sbs/sbs4/)*
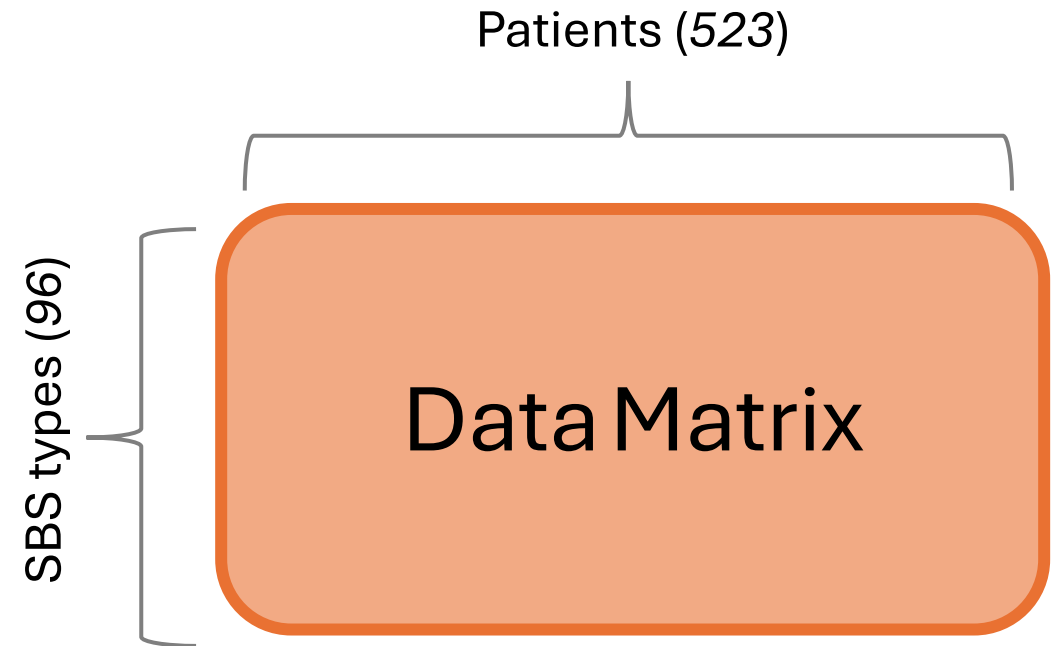
# Data

Mutational catalogue of **ovarian cancer** from the **Genomics England** 100,000 Genomes Project (GEL)

**523** whole genome sequences
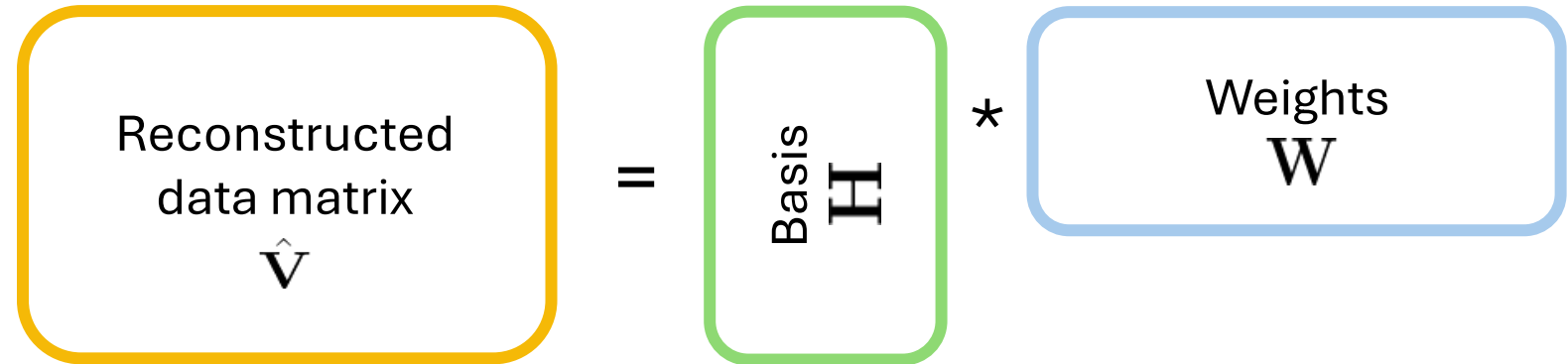
Single Base Substitutions (SBS)

**COSMIC** v 3.4 SBS GRCh37 for signature *comparison*

# Non-Negative Matrix Factorization (NMF)

$V \approx HW$

NMF

Reconstructed data matrix $\hat{\mathbf{V}}$ = Basis $\mathbf{H}$ * Weights $\mathbf{W}$

$V \approx VW_1W_2$

C-NMF/Autoencoder

Reconstructed data matrix $\hat{\mathbf{V}}$ = Data matrix $\mathbf{V}$ Encoder $\mathbf{W}_1$ * Decoder $\mathbf{W}_2$
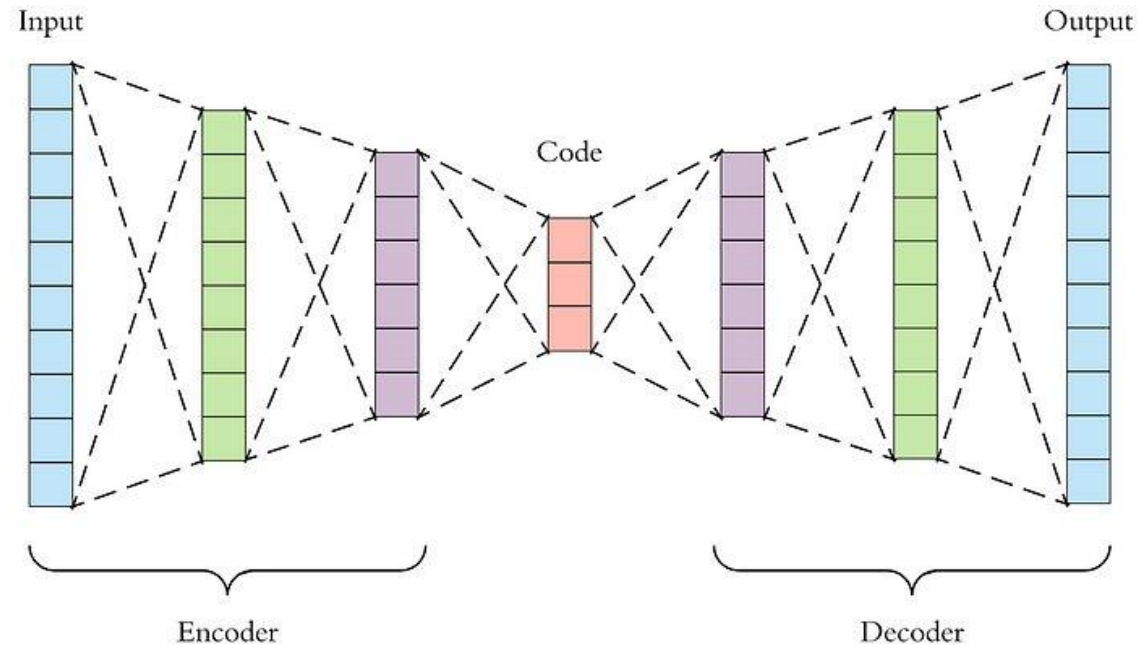
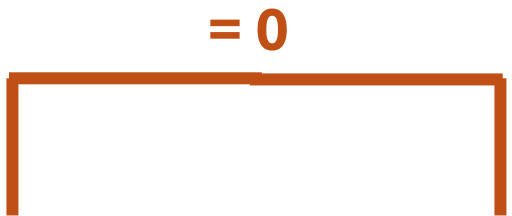$\mathbf{H}$ Basis

# Autoencoder

**Autoencoder**: neural network used for unsupervised learning to encode and decode data

- **Encoder:** compresses input into a lower-dimensional representation
- **Latent Space:** the compressed feature representation
- **Decoder:** reconstructs the input from the latent space

# Non-negative autoencoder & Convex NMF *equivalence*

**= 0**

**Shallow autoencoder:** $y_{pred} = \phi_{dec}(\phi_{enc}(\mathbf{V}\mathbf{W}_{enc} + \mathbf{b}_{enc})\mathbf{W}_{dec} + \mathbf{b}_{dec})$

$\phi_{enc}, \phi_{dec} : x \mapsto x$

Weights constrained to be non-negative

**Convex NMF:** $y_{pred} = \mathbf{V}\mathbf{W}_1\mathbf{W}_2$

# Shallow-AE & PCA

**PCA** (Principal Component Analysis): method that finds the best lower-dimensional representation of data while preserving as much variance as possible

A shallow autoencoder with a linear activation function (identity) is mathematically **equivalent** to PCA:

- *Linear transformations*: autoencoder with an identity activation function only learns linear transformation, like PCA

- *Eigenvectors & Weights*: encoder learns a weight matrix that spans the same subspace as the PCA principal components

- *Recostruction Error*: the objective of training the autoencoder with MSE is to minimize the reconstruction error, same as PCA
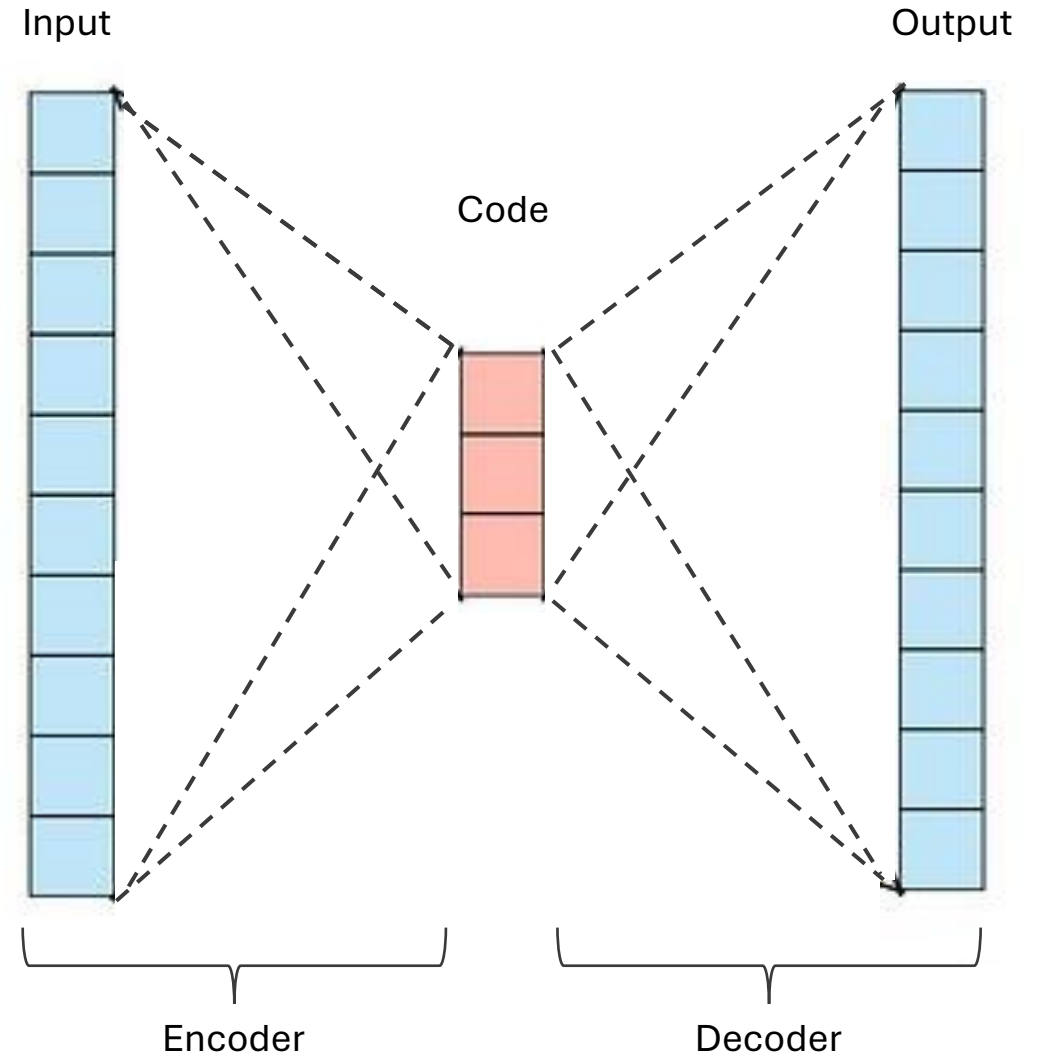
PCA:

$$min_W||X - XW^TW||^2$$

Shallow-AE:

$$min_W||X - y_{pred}||^2 = ||X - XW_{enc}W_{dec}||^2$$

# Autoencoder NMF (AE-NMF)

- **Input dimension:** 418

- **Latent dimension:** 4

- **Activation function:** Identity

- **Weight initialization:** $\mathcal{U}(0,1)$

# Training Procedure

- **Loss:** Frobenius Norm $\left\| X - \hat{X} \right\|_F$

- **Early stopping:** tolerance threshold (1e-10) on difference between subsequent loss values

- **ADAM** optimizer: learning rate fixed at 0.001

- **Weight clamping**: after gradient step and optimization, negative weights values are set to 0

# Results AE-NMF vs NMF

- **Runs:** 30

- **Train-Test split :** 80% - 20%

- **Train error:** $\left\| X - \hat{X} \right\|_F$

- **Test procedure**: *freeze S,* randomly intialize *E,* minimize $\left\| X - \hat{X} \right\|_F$

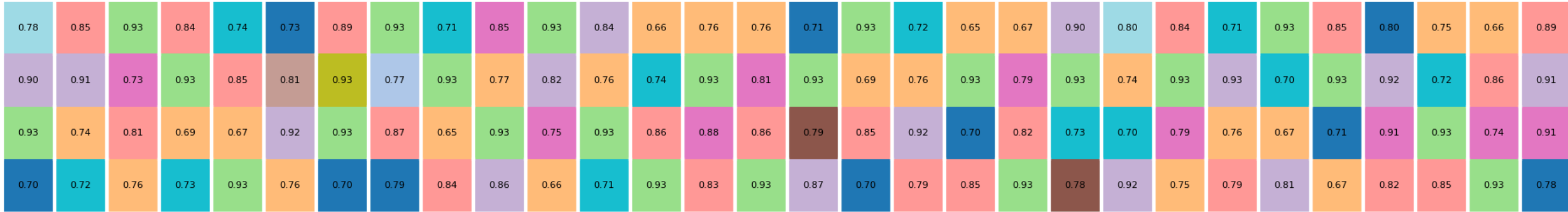| Reconstruction Error | NMF | AE-NMF |
|---|---|---|
| Average train losses | $1.48 * 10^4$ | $1.86 * 10^4$ |
| Average test losses | $1.05 * 10^5$ | $1.64 * 10^4$ |

# Cosine similarity with COSMIC

Reference for the extracted signature: **COSMIC v3.4 SBS GRCh37**

Cosine similarity matrix between *true* and *found* signatures

$$S_C(\tilde{\boldsymbol{h}}, \hat{\boldsymbol{h}}) = \frac{\tilde{\boldsymbol{h}} \cdot \hat{\boldsymbol{h}}}{||\tilde{\boldsymbol{h}}|| \, ||\hat{\boldsymbol{h}}||}$$

**Signature matching** is achieved through the **Hungarian algorithm** (linear assignment) ensuring the best available mapping
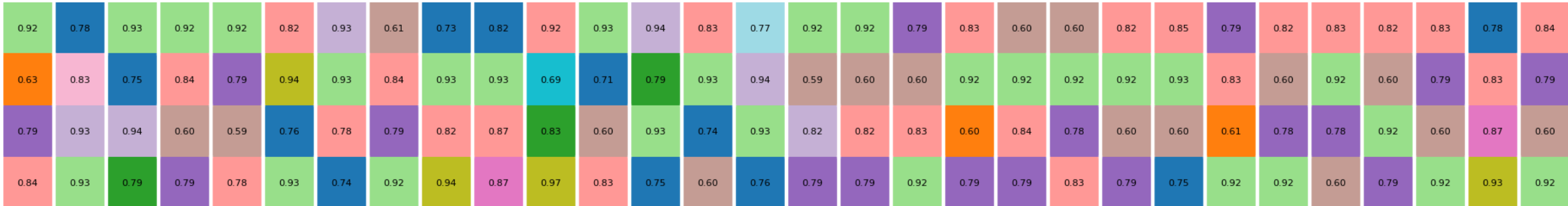
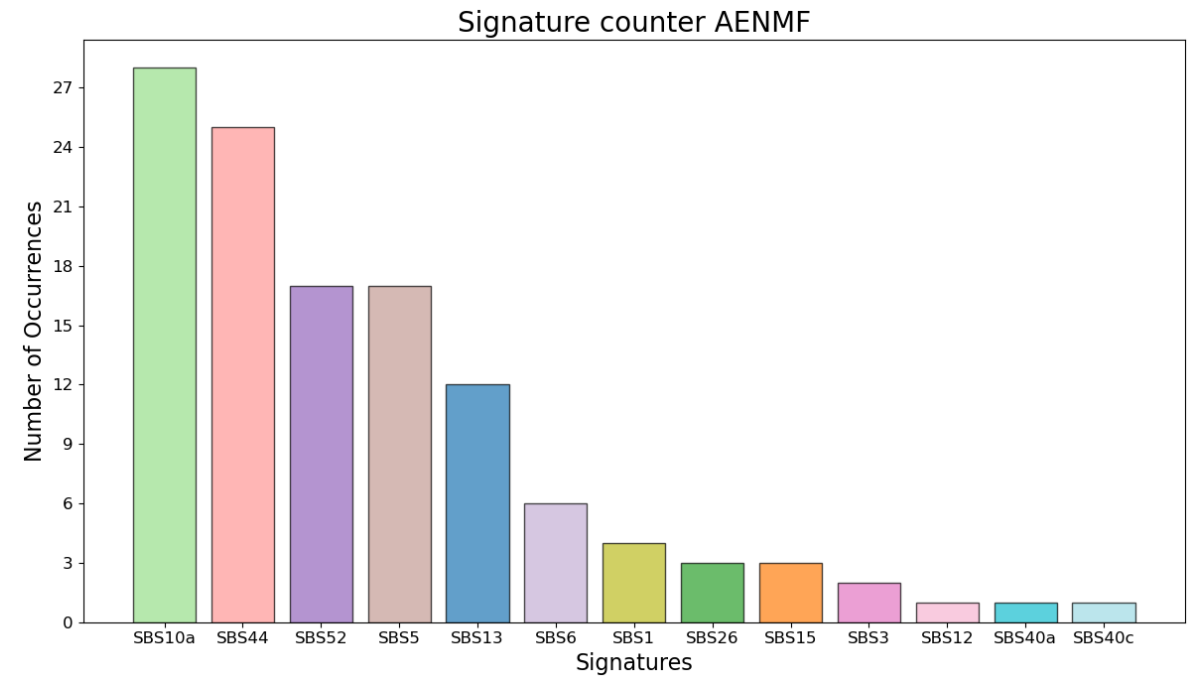# Cosine similarity comparison
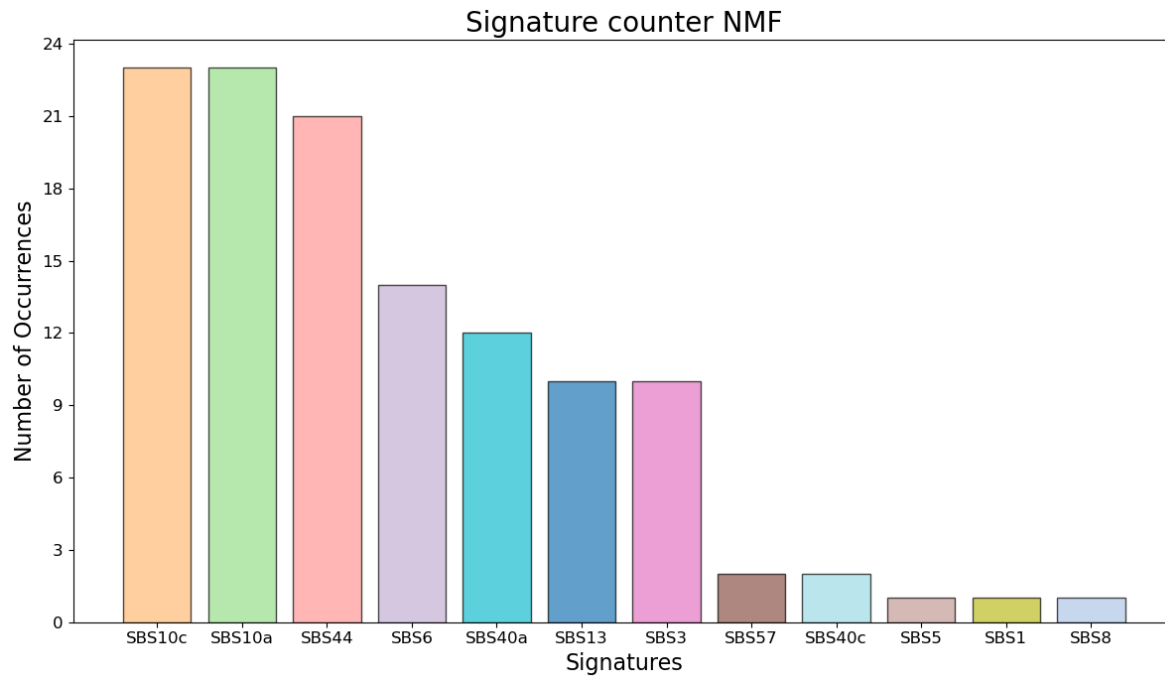


Cosine similarity matrix NMF

Matched signatures: SBS1, SBS10a, SBS10c, SBS13, SBS3, SBS40a, SBS40c, SBS44, SBS5, SBS57, SBS6, SBS8

Cosine similarity matrix AENMF

Matched signatures: SBS1, SBS10a, SBS12, SBS13, SBS15, SBS26, SBS3, SBS40a, SBS40c, SBS44, SBS5, SBS52, SBS6
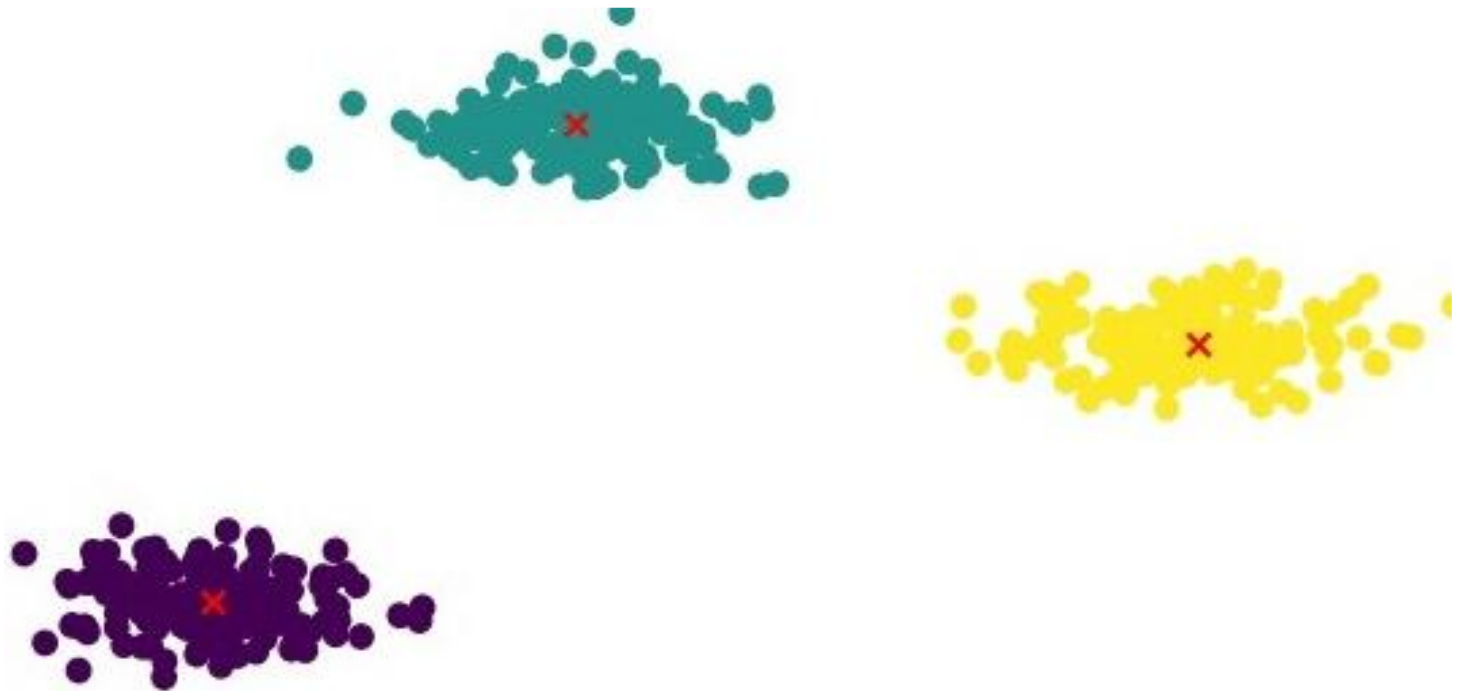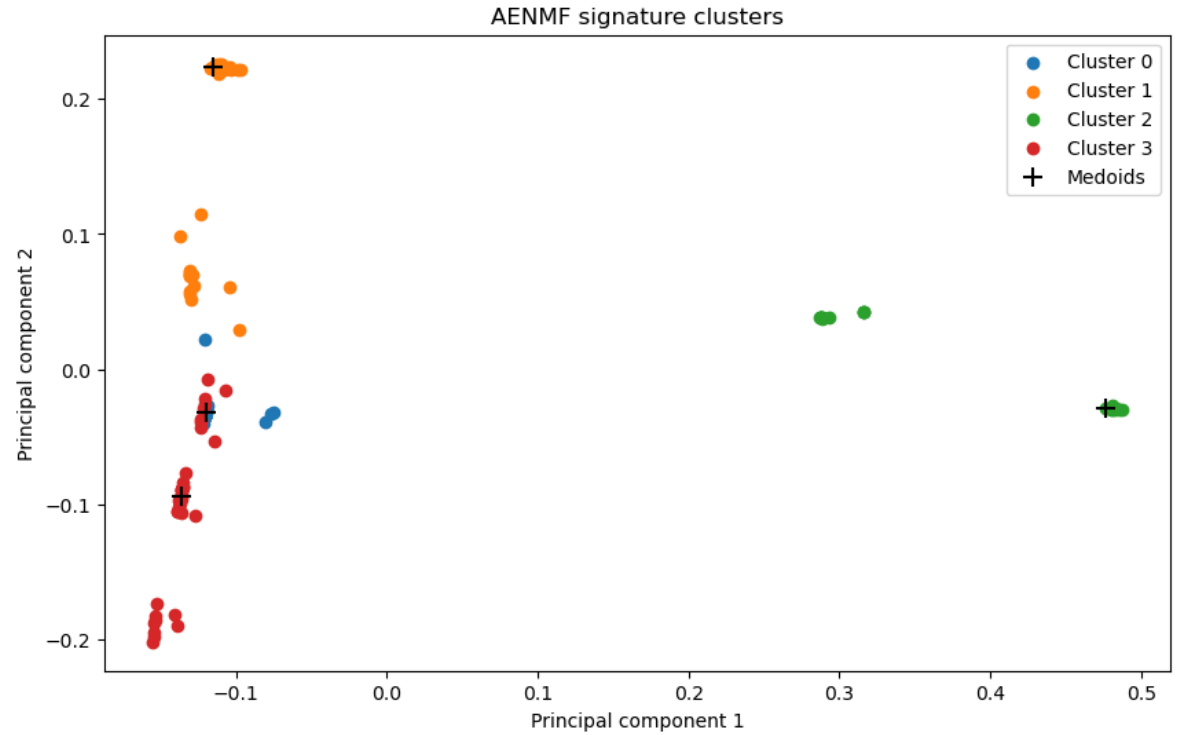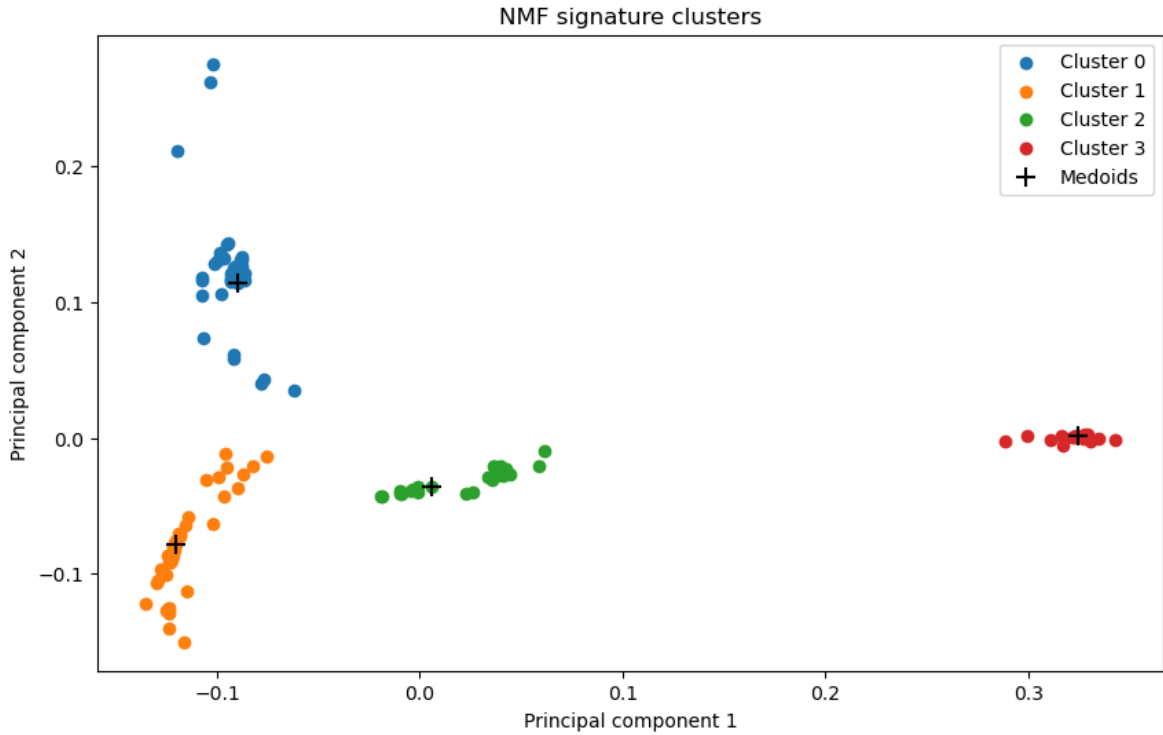
# Signature counts comparison

# Consensus signatures

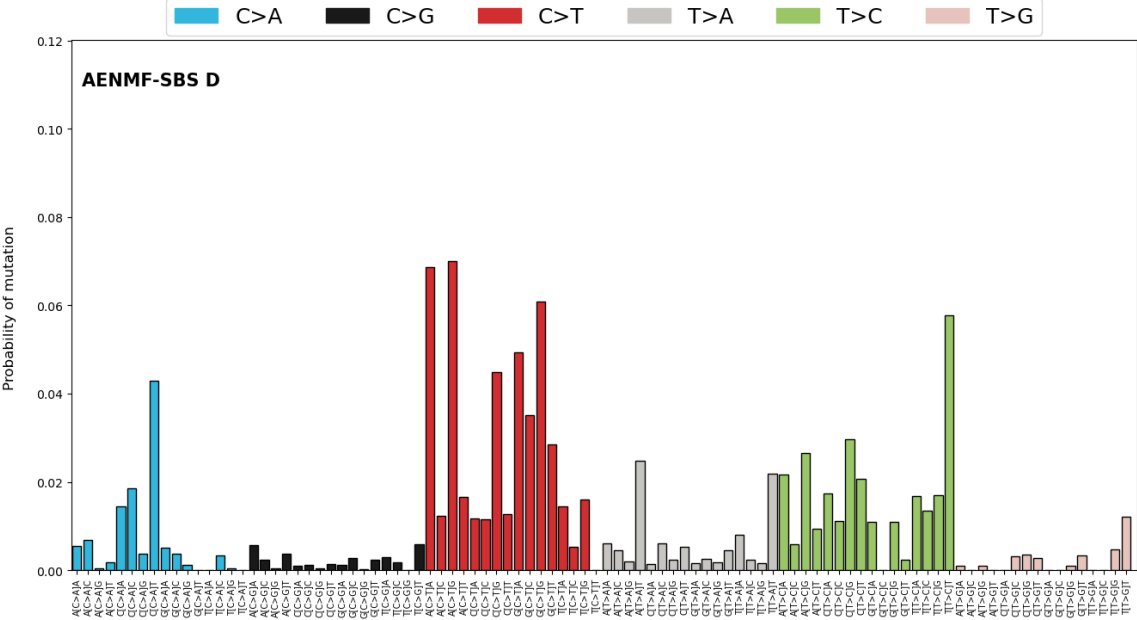Consensus signatures are found via a **K-medoids** algorithm based around the **cosine similarity** measure between the extracted signatures
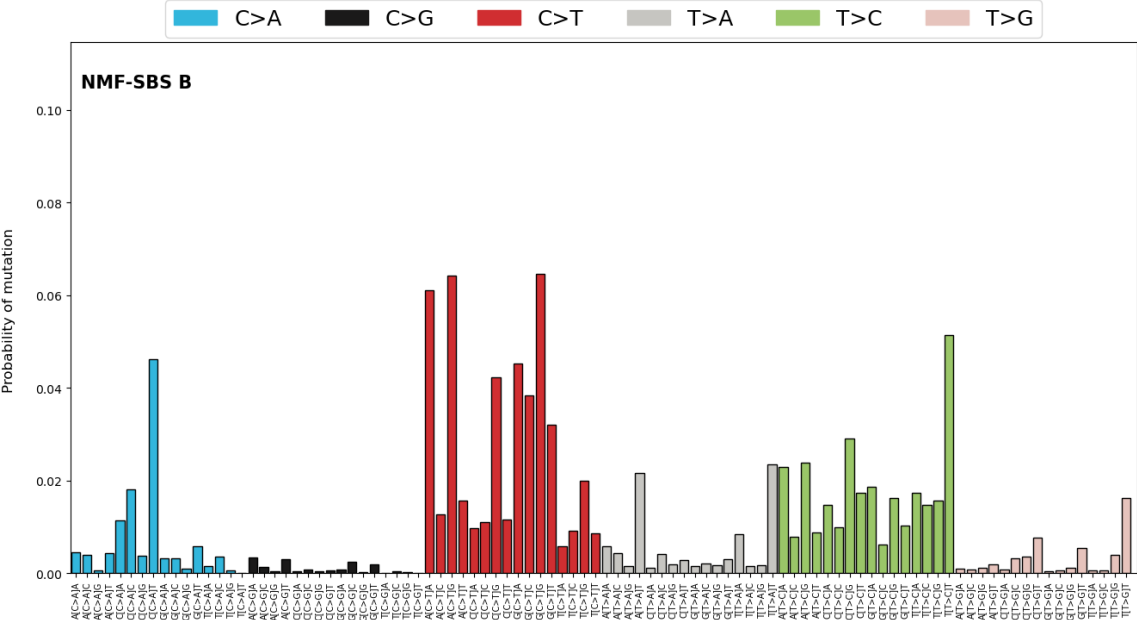
# Consensus signatures comparison

| NMF Extracted | True | Cosine Similarity |
|---|---|---|
| SBS-A | SBS40a | 0.74 |
| SBS-B | **SBS44** | **0.85** |
| SBS-C | SBS10c | 0.69 |
| SBS-D | **SBS10a** | **0.93** |

| AENMF Extracted | True | Cosine Similarity |
|---|---|---|
| SBS-A | SBS5 | 0.60 |
| SBS-B | SBS52 | 0.78 |
| SBS-C | **SBS10a** | **0.92** |
| SBS-D | **SBS44** | **0.83** |

# Comparison with COSMIC signature

# Non-linear autoencoder

## Why?

- A **non-linear autoencoder** can effectively model complex genomic interactions that could be missed by simpler, linear methods.

- A **Poisson-based** loss function naturally accommodates the count-based nature of mutational data.

- **Sparsity constraints** may improve the interpretability of extracted signatures, reducing overlap and facilitating biological interpretation.

# Input data adjustment

$m = 96$

$n = 523$

Data matrix
**V**

$$V \in \mathbb{R}^{m \times n}$$

$$V \approx HW$$
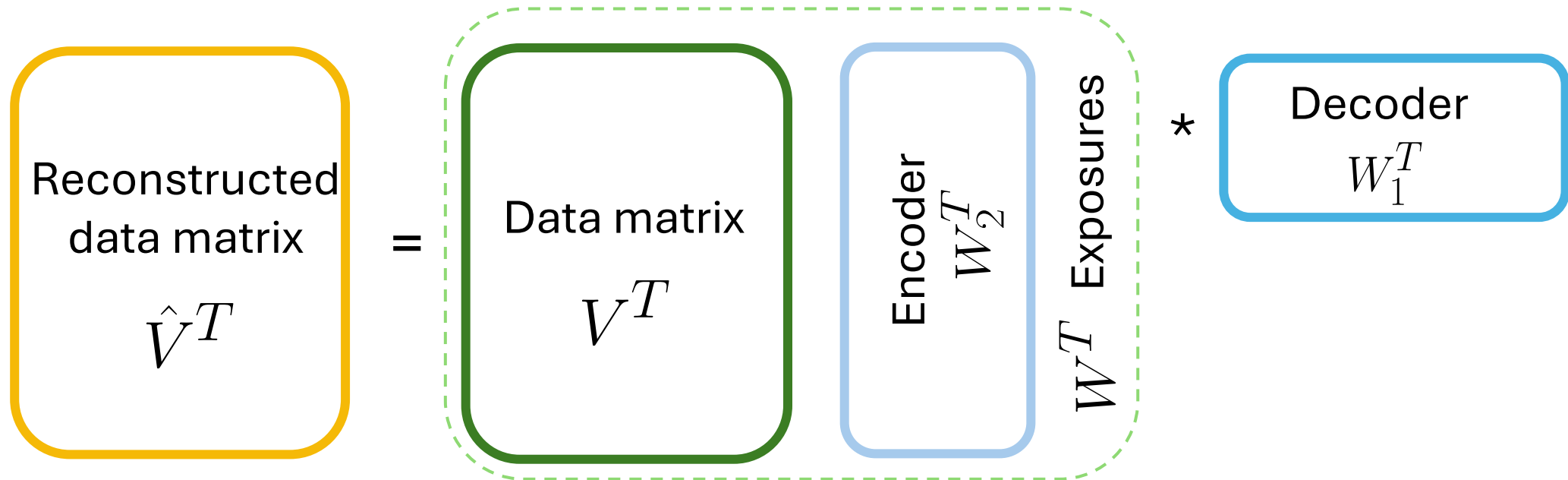
Transposed
data matrix
$V^T$

$$V^T \in \mathbb{R}^{n \times m}$$

$$V^T \approx W^T H^T$$

# Signature and exposure inversion



Reconstructed data matrix $\hat{V}^T$ = Data matrix $V^T$ | Encoder $W_2^T$ | Exposures $W^T$ * Decoder $W_1^T$

# Autoencoder Architecture

**Encoder**

**Input dimension:** 96

**Three** hidden layers: 128, 64, 32

**Latent dimension:** between 3 and 9

**Decoder (shallow)**

**One** linear layer

**Output dimension:** 96

# Activation function

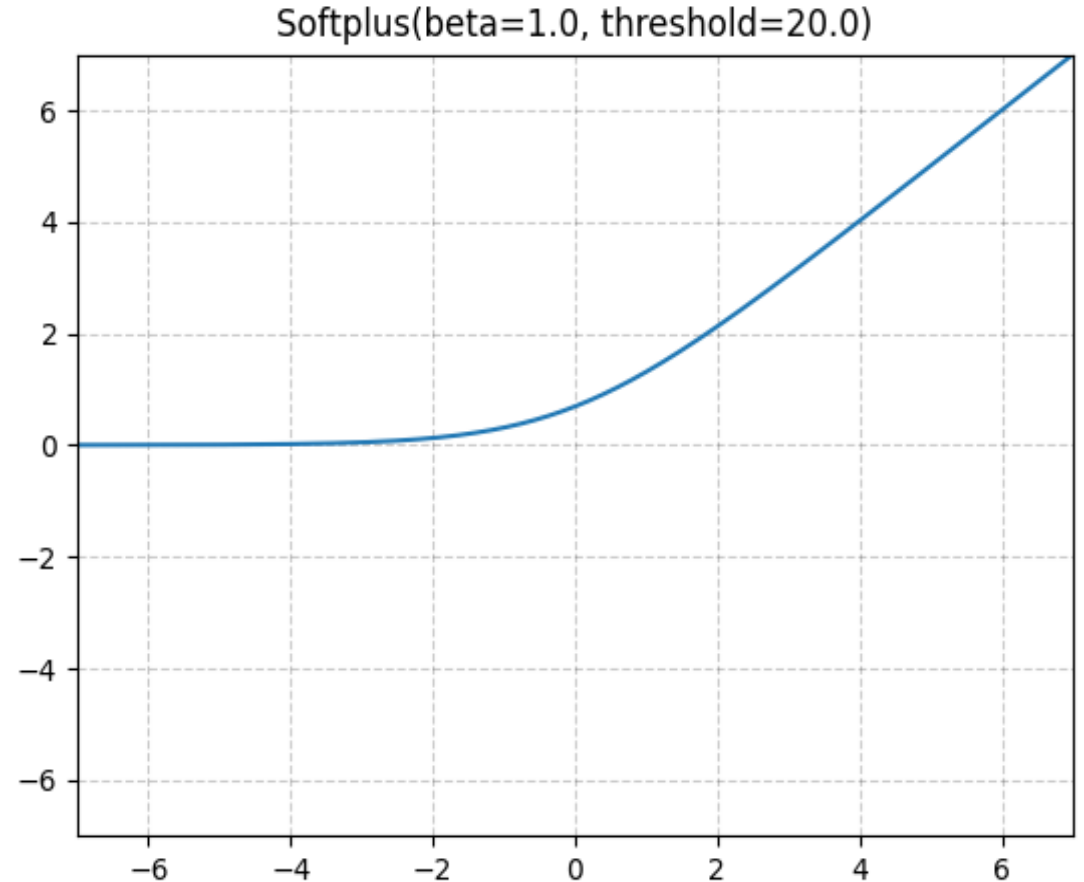The **Softplus** activation function was used in all encoder layers *and* in the *latent layer*.

For the decoder an **identity** activation function was used

Softplus(beta=1.0, threshold=20.0)

$$\frac{1}{\beta} \log \left(1 + \exp(x \cdot \beta)\right)$$

# Weight Initialization

The *xavier uniform* method was used to initialize all the layers weights

Clamping on the *decoder* weights to avoid negative values

$$w_{ij} \sim U\left(-\sqrt{\frac{6}{n_{\text{in}} + n_{\text{out}}}}, \sqrt{\frac{6}{n_{\text{in}} + n_{\text{out}}}}\right)$$

# Loss optimization

Non-negative Poisson Likelihood

Signature must be non negative

$$L(x; \hat{x}) = -x \log(\hat{x}) + \hat{x} + \beta \log \left( \det(WW^\top + I) \right) \quad \text{subject to} \quad W \geq 0$$

Minimum Volume Regularizer

# Multinomial Bootstrapping

Each i-th patient of the **augmented dataset** is drawn from the
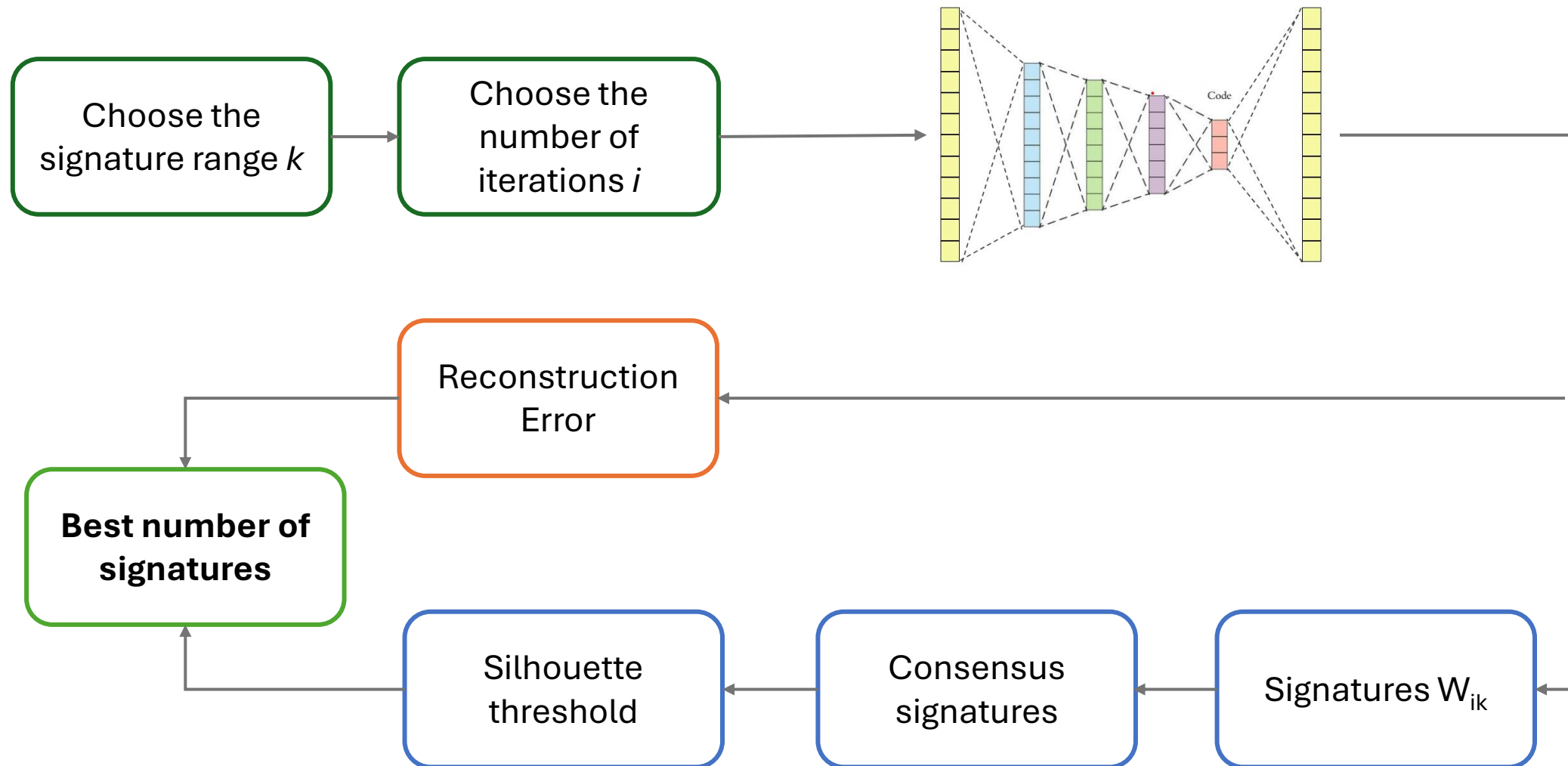*i-th* Multinomial distribution $S_i \sim M_i(N_i, p_i)$ parametrized by:

$N_i$    Total **number** of mutations for patient *i* in the original dataset

$p_i$    **Probablity** of each mutation for patient *i* in the original dataset

We generated **50 new datasets** and stacked them to create one larger augmented dataset to train the autoencoder with.

*The original dataset is used as a validation set* to check for overfitting via earlystopping

# Choosing the correct number of signatures

# Choosing the correct number of signatures

**Consensus signatures**

For each group of *n* mutational signatures run the *K-medoids* algorithm to find the consensus signatures



**Silhouette threshold**

Discard the consensus signatures which minimum and average silhouette score don't satisfy the threshold (0.2 - 0.5)

$$th_{min} \geq \rho_{min}$$
$$th_{avg} \geq \rho_{avg}$$

**Reconstruction Error**

Computed using the *Frobenius* norm on the difference between the original data and the reconstructed one. Used for early stopping via a patience counter

$$\left\| X - \hat{X} \right\|_F$$

# Bayesian Optimization

We want to efficiently explore the hyperparameter space towards the global optimum.

Goal: find the set of input parameters $x$ that **maximize** a function $f(x)$

$$\underset{x \in A}{\mathrm{argmax}}\, f(x)$$

Components:

- **Surrogate model**: a statistical model to **approximate** the objective function $f(x)$ → *Gaussian Processes*

- **Acquisition function**: function that **guides** where to sample next → *Expected Improvement* (expected value of how much better the function value $f(x)$ at a given point $x$ is compared to the best known function value achieved at step $n$)

# Hyperparameter Tuning

| Hyperparameters | Range |
|---|---|
| ADAM learning rate $lr$ | [1e-4, 1e-2] |
| Regularization strength $\beta$ | [1e-4, 1e-2] |

# Pipeline

# Results

| | Extracted | True | Cosine Similarity |
|---|---|---|---|
| $k = 5$ | SBS-A | SBS10a | 0.85 |
| | SBS-B | SBS40a | 0.87 |
| lr = 0.008 | SBS-C | SBS10c | 0.75 |
| | SBS-D | SBS56 | 0.89 |
| B = 0.005 | SBS-E | SBS6 | 0.82 |

# Cosine similarity comparison



Cosine similarity matrix non-linear AE

# Signatures comparison

# Comparison with COSMIC signature



**Final reconstruction error**

$4.3 * 10^4$

# Aetiology

| Signature | Times found | Aetiology |
|---|---|---|
| **SBS10a** | 29/30 | Polymerase epsilon (POLE) exonuclease domain mutations |
| **SBS10c** | 24/30 | POLE exonuclease domain mutations |
| **SBS56** | 22/30 | Possible sequencing artifacts |
| **SBS3** | 20/30 | Defective homologous recombination-based DNA (BRCA1 and BRCA2) |
| **SBS6** | 17/30 | Defective DNA mismatch repair |
| **SBS40a** | 10/30 | Unknown |

# Conclusions

Overall, NLAE demonstrated consistency in its findings, repeatedly extracting the same signatures across multiple runs. It shares many signatures with NMF and may have captured an acquisition error.

The reconstruction error increased slightly, but it remains within the same order of magnitude as other methods.

The experiment was generally successful, but further validation on diverse datasets is needed to strengthen our conclusions. Additionally, exploring other techniques to mitigate overfitting could provide valuable insights.

Finally, investigating overdispersion in the reconstruction error and eventually adjusting the Poisson term in the loss to be a Negative Binomial one could be of interest

# THANK YOU !

# References:

- On the Relation Between Autoencoders and Non-negative Matrix Factorization, and Their Application for Mutational Signature Extraction ; Egendal et al. (2024)

- MUSE-XAE: MUtational Signature Extraction with eXplainable AutoEncoder enhances tumour types classification; Pancotti et al. (2024)

- Decoding whole-genome mutational signatures in 37 human pan-cancers by denoising sparse autoencoder neural network; Pei et al. (2020)