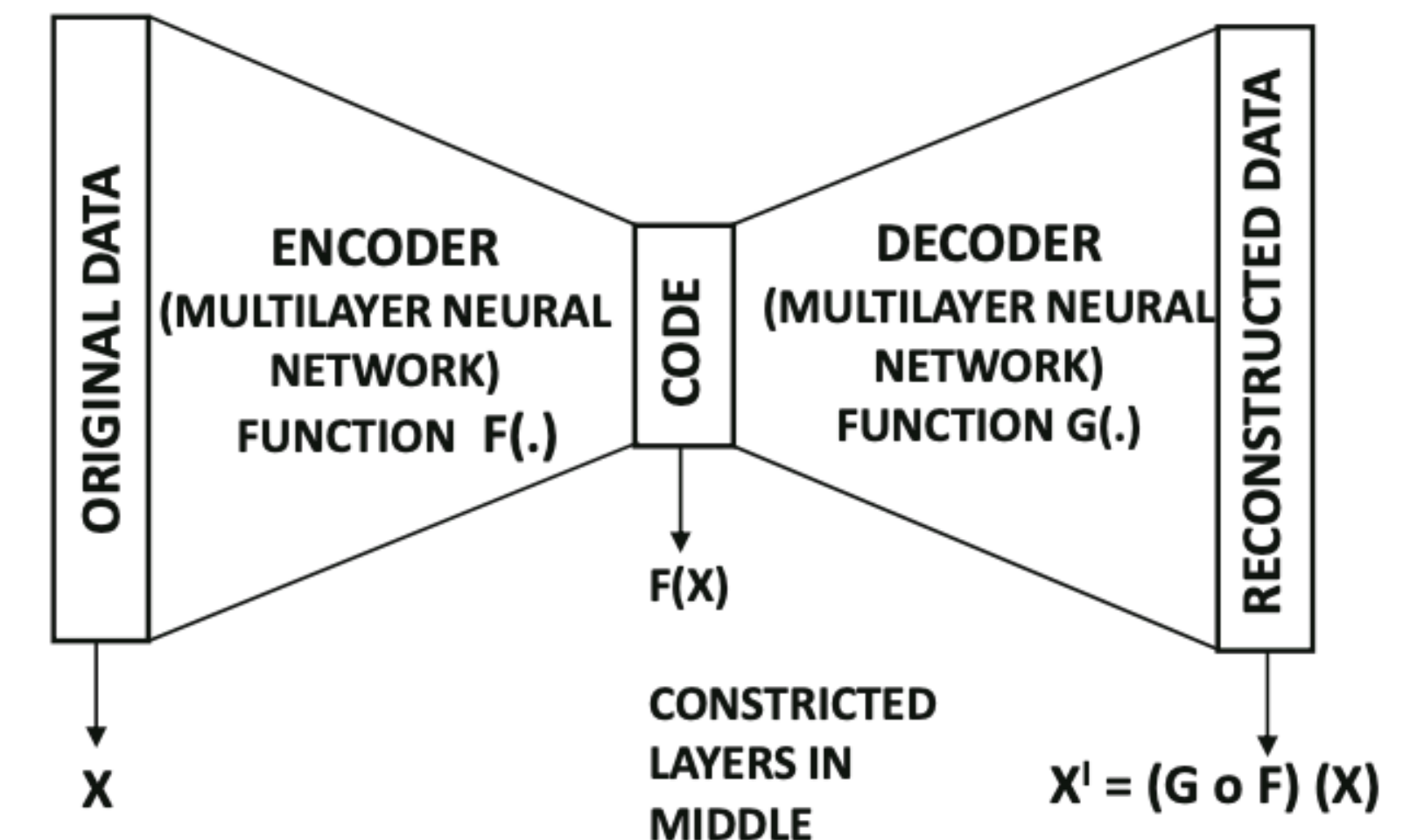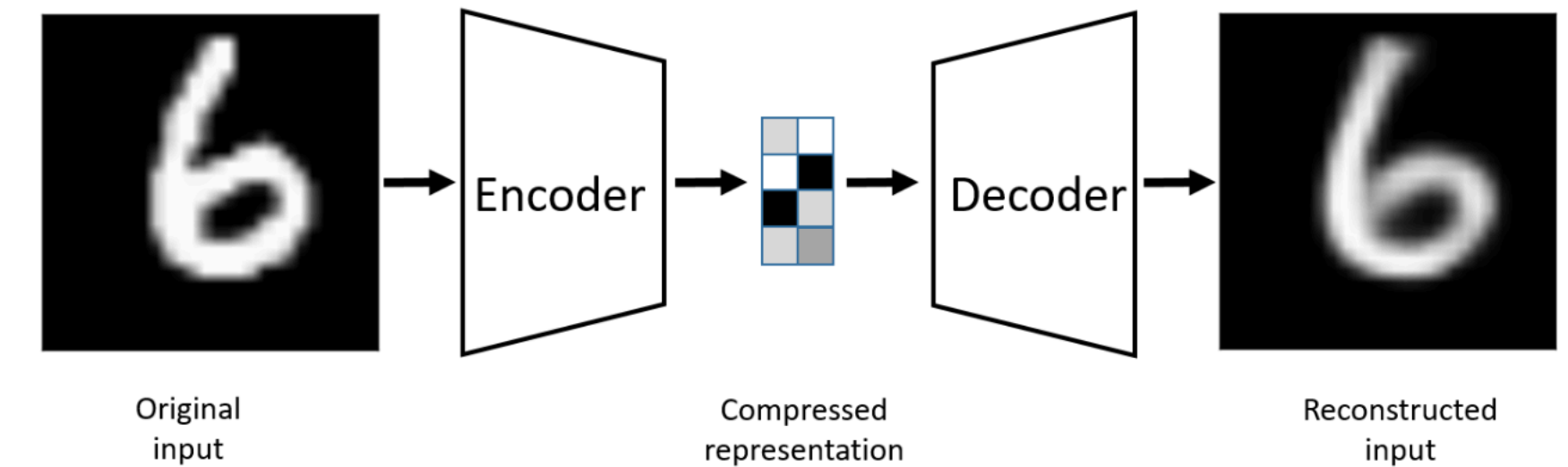19-11–2024

# Projects

- Can be done in a group (max two students)

- Be careful about your project partner!

- If he is auditing the course then you will be in trouble!

- Define your own project

- Submit a one page project proposal- within fixed time (first four weeks)?

- Finished the work within the time-line

- Report submission
  - Submission deadline: <span style="color:red">seven days before the final exam date,</span> is strict and you can adjust your assignment buffer days here - <span style="color:red">24-11-2024</span>
  - We will consider 11:59PM as our day end

- Final presentation
  - 20 min (divided into group members)
  - <span style="color:red">Five days before the final exam date - 26-11-2024 & 27-11-2024</span>

# Auto-encoder

- Dimensionality reduction:
  - Classical method (linear): PCA
  - Data matrix: $X \in \mathbb{R}^{d \times N}$
  - Projection matrix: $W \in \mathbb{R}^{d \times d_1}$
  - Lower dims. representation (LDR): $Z = W^T X; Z \in \mathbb{R}^{d_1 \times N}$
  - Reconstruction: $\bar{X} = [WW^T]^{-1} WZ$
  - In general - $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d_1}; g : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^d$
  - LDR: $Z = f(X)$
  - Reconstruction: $\bar{X} = g(Z) = g[f(X)]$
  - Non-linear dimensionality reduction



Images: Bank et al. and Chart Agarwal

# Auto-encoder (cont…)

- Issues in Auto-encoder - $f : \mathbb{R}^d \to \mathbb{R}^{d_1}; g : \mathbb{R}^{d_1} \to \mathbb{R}^d$

  - Overfitting

    - What about the size of the latent dims. $(d_1)$ ?

    - $d < d_1$

    - $d = d_1$

    - $d > d_1$

    - $d_1 = 1$

  - Bias-variance tradeoff

    - We want the architecture of the auto-encoder to be able to reconstruct the input well

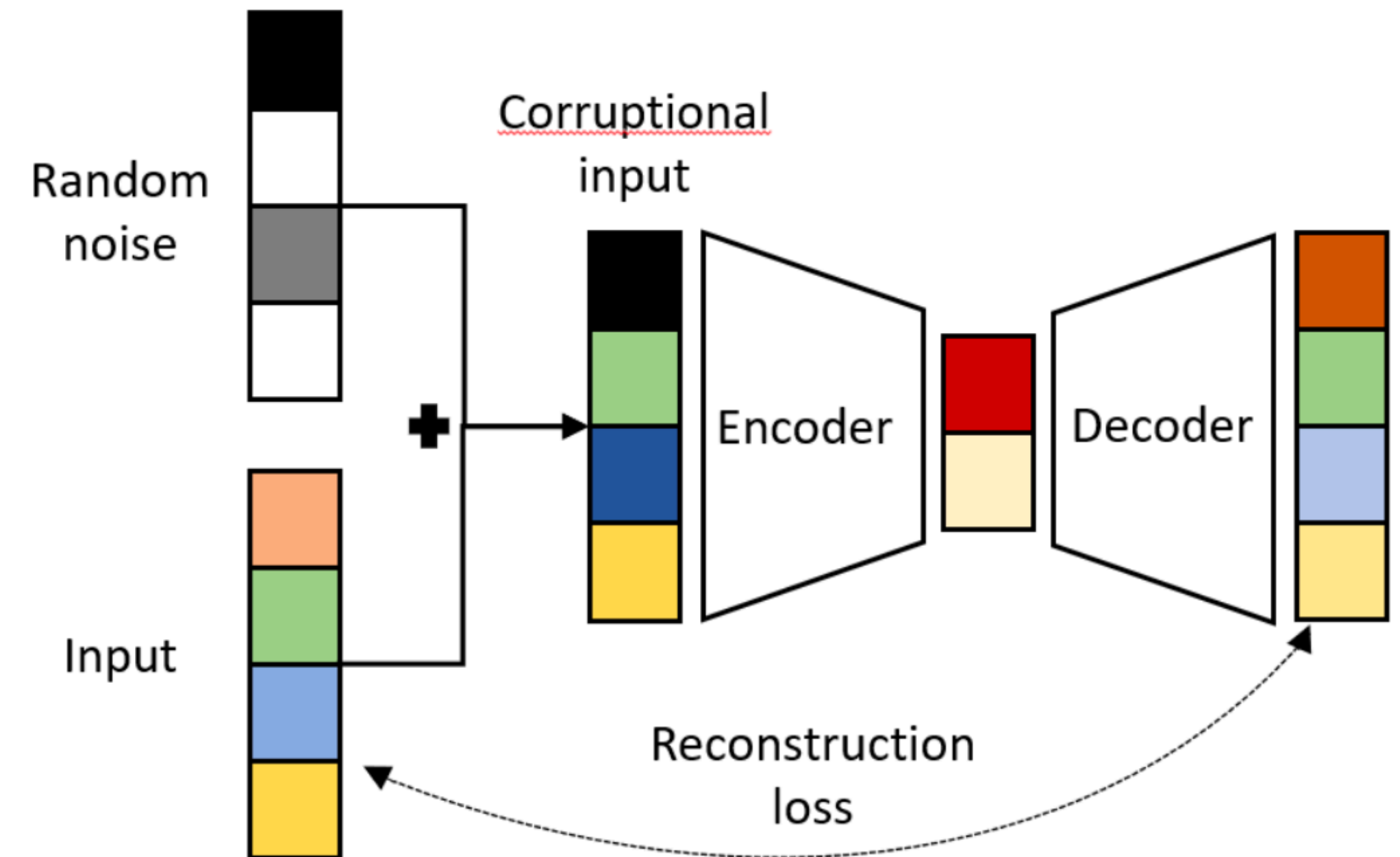    - We want the low representation to generalise to a meaningful one

  - Solution:

    - Sparse Auto-encoder; Denoising Auto-encoder; …



Original input    Compressed representation    Reconstructed input

# Auto-encoder (cont...)

- Denoising Auto-encoder:
  - Add small noise to the input data
  - What about the random noise ?
    - $\tilde{X} = X + \epsilon; \epsilon \sim \mathcal{N}(0,\sigma^2)$
    - $C_\sigma(\tilde{X}|X) = \mathcal{N}(X, \sigma^2 I)$
    - $C_p(\tilde{X}|X) = \beta \odot X; \beta \sim Ber(p)$

# Auto-encoder (cont…)

- Demo:
  - ‣ https://cs.stanford.edu/~karpathy/convnetjs/demo/autoencoder.html

# Auto-encoder (cont…)

- Take-home on Auto-encoder
  - ‣ Non-linear dimensionality reduction
  - ‣ Useful for unsupervised feature learning



Original input     Encoder     Compressed representation     Decoder     Reconstructed input

# Variational auto-encoder



$$logp_\theta(x) = E_z[p_\theta(x|z)] - KL[q_\phi(z|x)\|p_\theta(z)] + \boxed{\begin{array}{c} KL[q_\phi(z|x)\|p_\theta(z|x)] \\ \geq 0 \end{array}}$$

$$= E_z[p_\theta(x|z)] - KL[q_\phi(z|x)\|p_\theta(z)]$$

# Variational auto-encoder

- Latent space

# Molecule representation in continuous space

[1]**Samanta** et al. VAE-Sim: A Novel Molecular Similarity Measure Based on a Variational Autoencoder. *Molecules* 2020, *25*, 3446
[2]Yash Khemchandani, Steve O'Hagan, **Soumitra Samanta**, Neil Swainston, Timothy J. Roberts, Danushka Bollegala and Douglas B. Kell, **DeepGraphMolGen, a multi-objective, computational strategy for generating molecules with desirable properties: a graph convolution and reinforcement learning approach**, In Journal of Cheminformatics, 12, 53, September 2020

# Generative model for new molecule generation: Variational Autoencoder

[1]Samanta et al. VAE-Sim: A Novel Molecular Similarity Measure Based on a Variational Autoencoder. Molecules 2020, 25, 3446
[2]Yash Khemchandani, Steve O'Hagan, Soumitra Samanta, Neil Swainston, Timothy J. Roberts, Danushka Bollegala and Douglas B. Kell, DeepGraphMolGen, a multi-objective, computational strategy for generating molecules with desirable properties: a graph convolution and reinforcement learning approach, In Journal of Cheminformatics, 12, 53, September 2020

# Results

- ZINC15: 2D Drug-Like, clean and in-stock (6,202,415 substances)
- Data partition: randomly partition the dataset into train-50% (3,101,207), validation-20% (1,240,483) and test-30% (1,860,725)

| Data partition | #Samples | #Valid reconstructed samples | Accuracy(%) |
|----------------|----------|------------------------------|-------------|
| Train | 3,101,207 | 2,984,669 | 96.24 |
| Validation | 1,240,483 | 1,180,189 | 95.13 |
| Test | 1,860,725 | 1,771,064 | 95.18 |

# Reconstructed molecules on test dataset



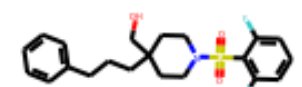mol_id_0_original    mol_id_0_reconstructed    mol_id_1_original    mol_id_1_reconstructed    mol_id_2_original    mol_id_2_reconstructed    mol_id_3_original    mol_id_3_reconstructed    mol_id_4_original    mol_id_4_reconstructed
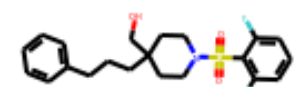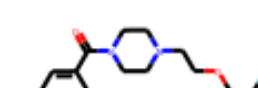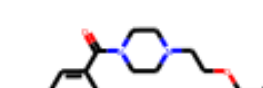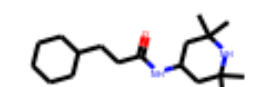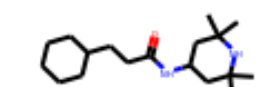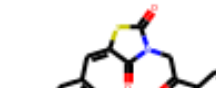
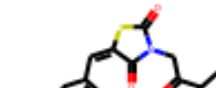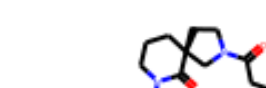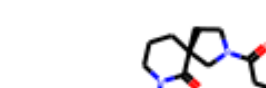mol_id_5_original    mol_id_5_reconstructed    mol_id_6_original    mol_id_6_reconstructed    mol_id_7_original    mol_id_7_reconstructed    mol_id_8_original    mol_id_8_reconstructed    mol_id_9_original    mol_id_9_reconstructed
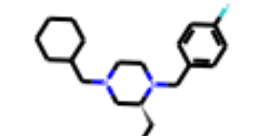
mol_id_10_original    mol_id_10_reconstructed    mol_id_11_original    mol_id_11_reconstructed    mol_id_12_original    mol_id_12_reconstructed    mol_id_13_original    mol_id_13_reconstructed    mol_id_14_original    mol_id_14_reconstructed

mol_id_15_original    mol_id_15_reconstructed    mol_id_16_original    mol_id_16_reconstructed    mol_id_17_original    mol_id_17_reconstructed    mol_id_18_original    mol_id_18_reconstructed    mol_id_19_original    mol_id_19_reconstructed
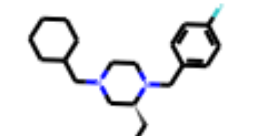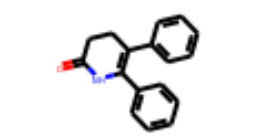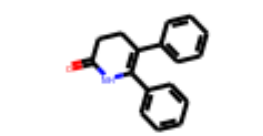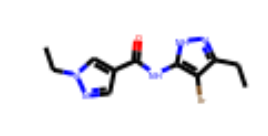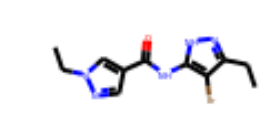
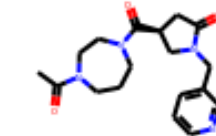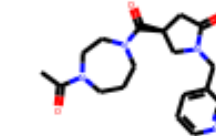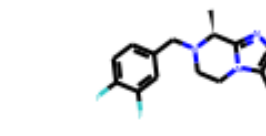mol_id_20_original    mol_id_20_reconstructed    mol_id_21_original    mol_id_21_reconstructed    mol_id_22_original    mol_id_22_reconstructed    mol_id_23_original    mol_id_23_reconstructed    mol_id_24_original    mol_id_24_reconstructed
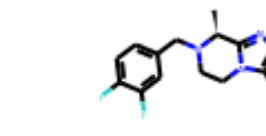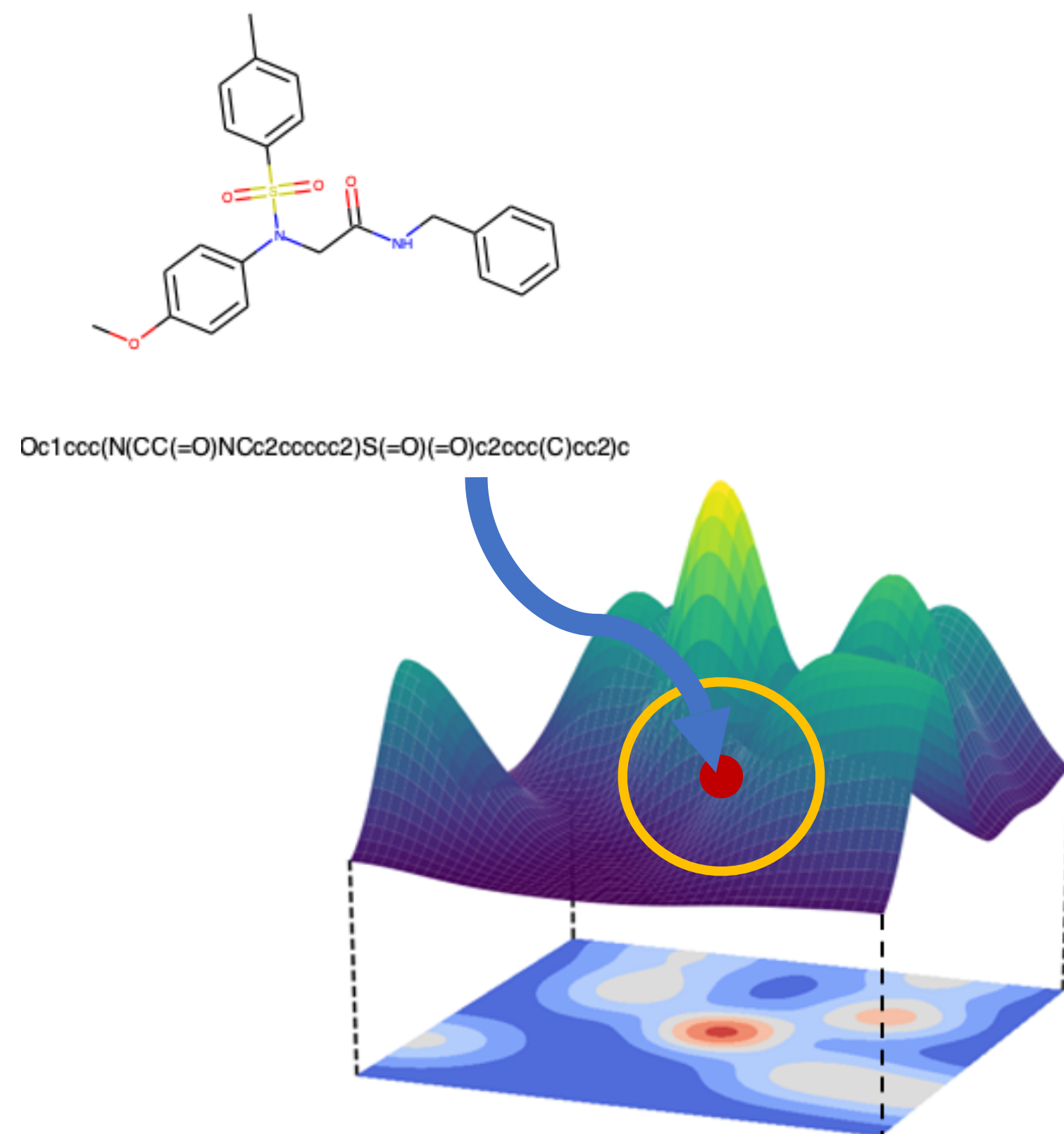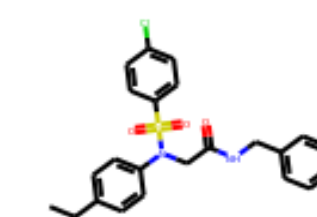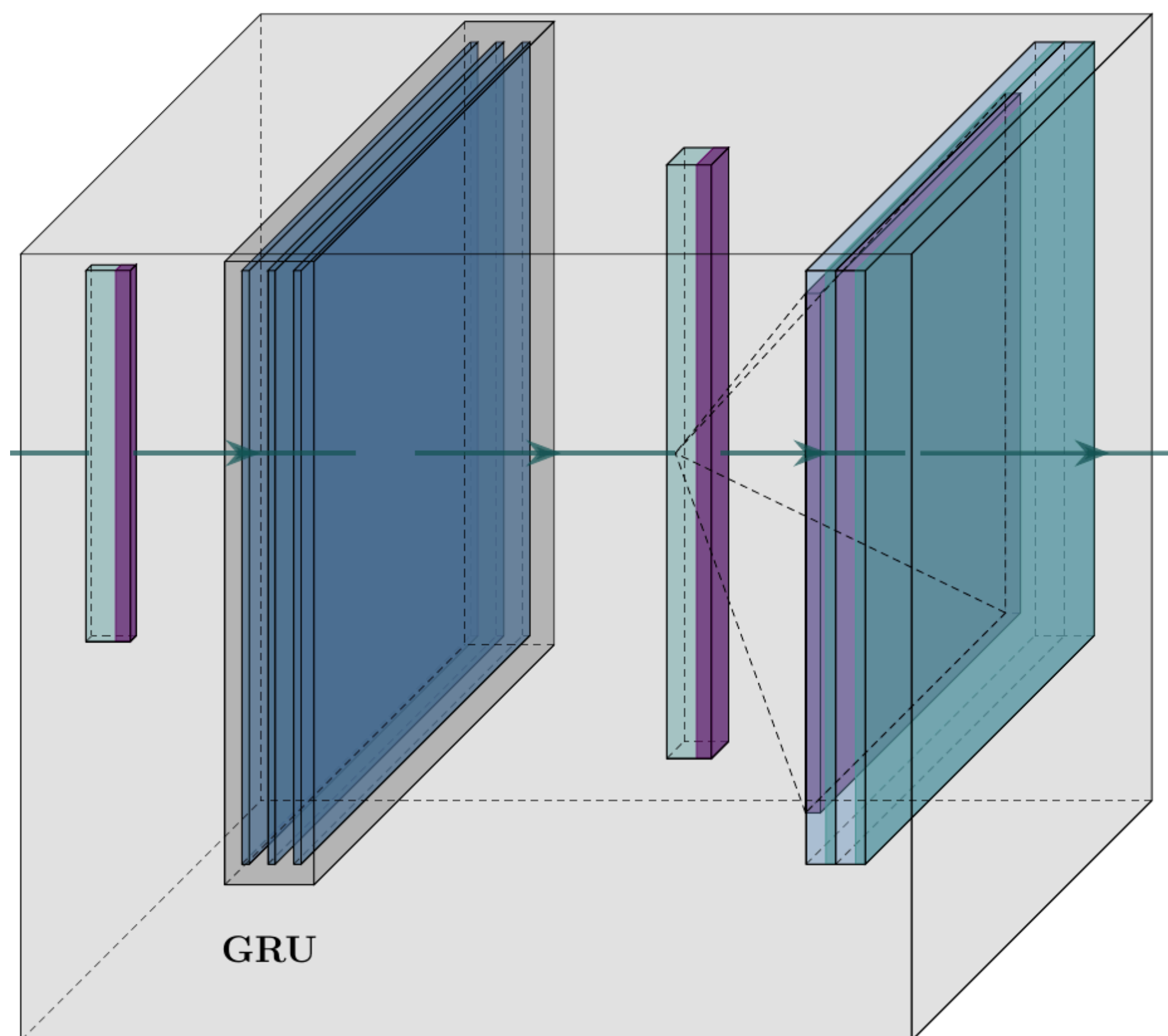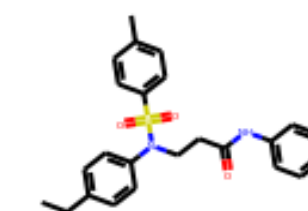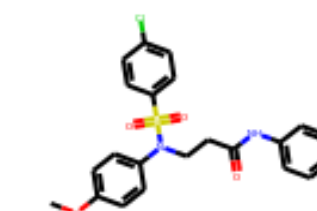
# Sampling molecules from the latent space

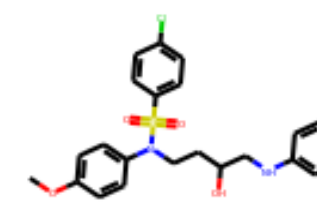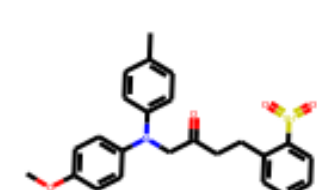Oc1ccc(N(CC(=O)NCc2ccccc2)S(=O)(=O)c2ccc(C)cc2)c

Sample from
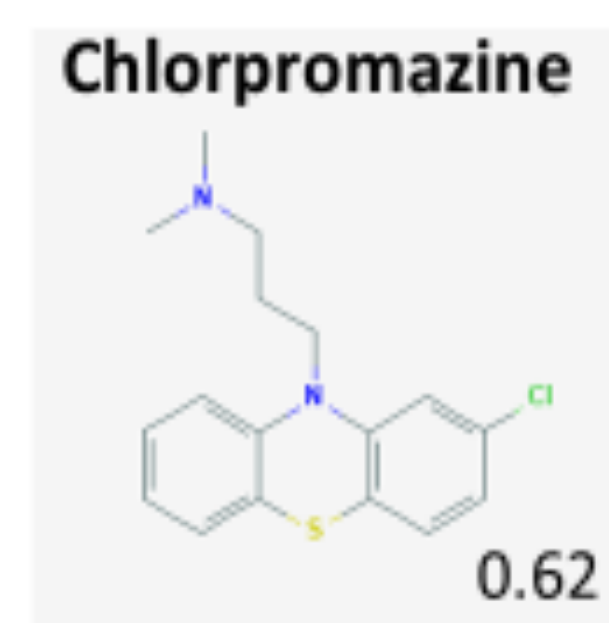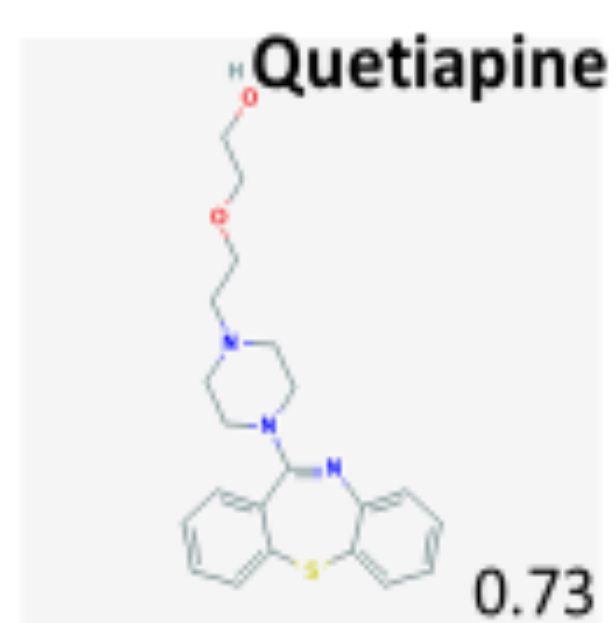latent space

GRU

0.06584072

0.11507881

0.12638378

0.13463712

0.16970044

Output

# Nearest neighbour sampling for a test data



| q_molecule | 0.06584072 | 0.06959683 | 0.070566654 | 0.07254505 | 0.074860275 | 0.07577288 | 0.07663196 | 0.10096228 | 0.10829282 |

| 0.11299068 | 0.11507881 | 0.116974175 | 0.12149018 | 0.12172425 | 0.122670054 | 0.123652816 | 0.12388712 | 0.12460661 | 0.12626642 |

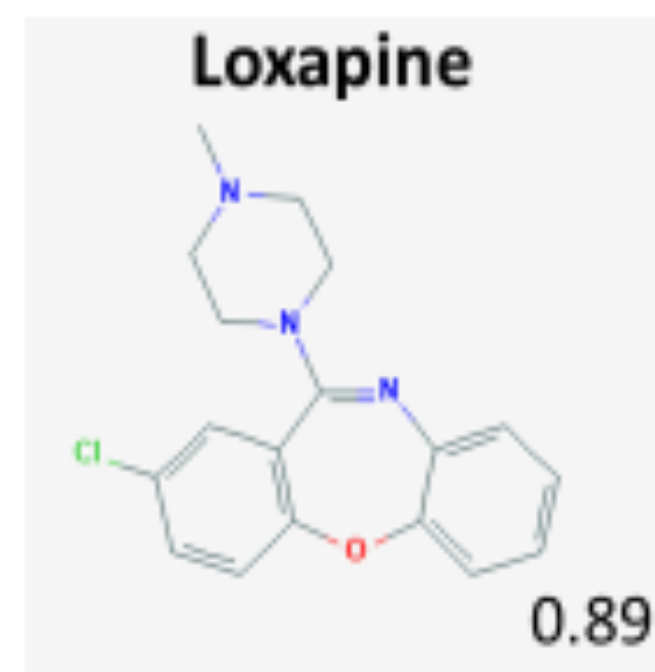| 0.12626886 | 0.12638378 | 0.12666166 | 0.12687105 | 0.1283449 | 0.12876374 | 0.12989932 | 0.13037455 | 0.13097942 | 0.13107568 |

| 0.1315139 | 0.13463712 | 0.14074236 | 0.1412546 | 0.14372319 | 0.15435714 | 0.16187131 | 0.16263306 | 0.16659111 | 0.1675002 |

| 0.16769993 | 0.16970044 | 0.17048365 | 0.17070031 | 0.17364568 | 0.17403269 | 0.17471558 | 0.1752078 | 0.17576742 | 0.17604893 |

# Nearest neighbour search

# Latent dimension evaluation