

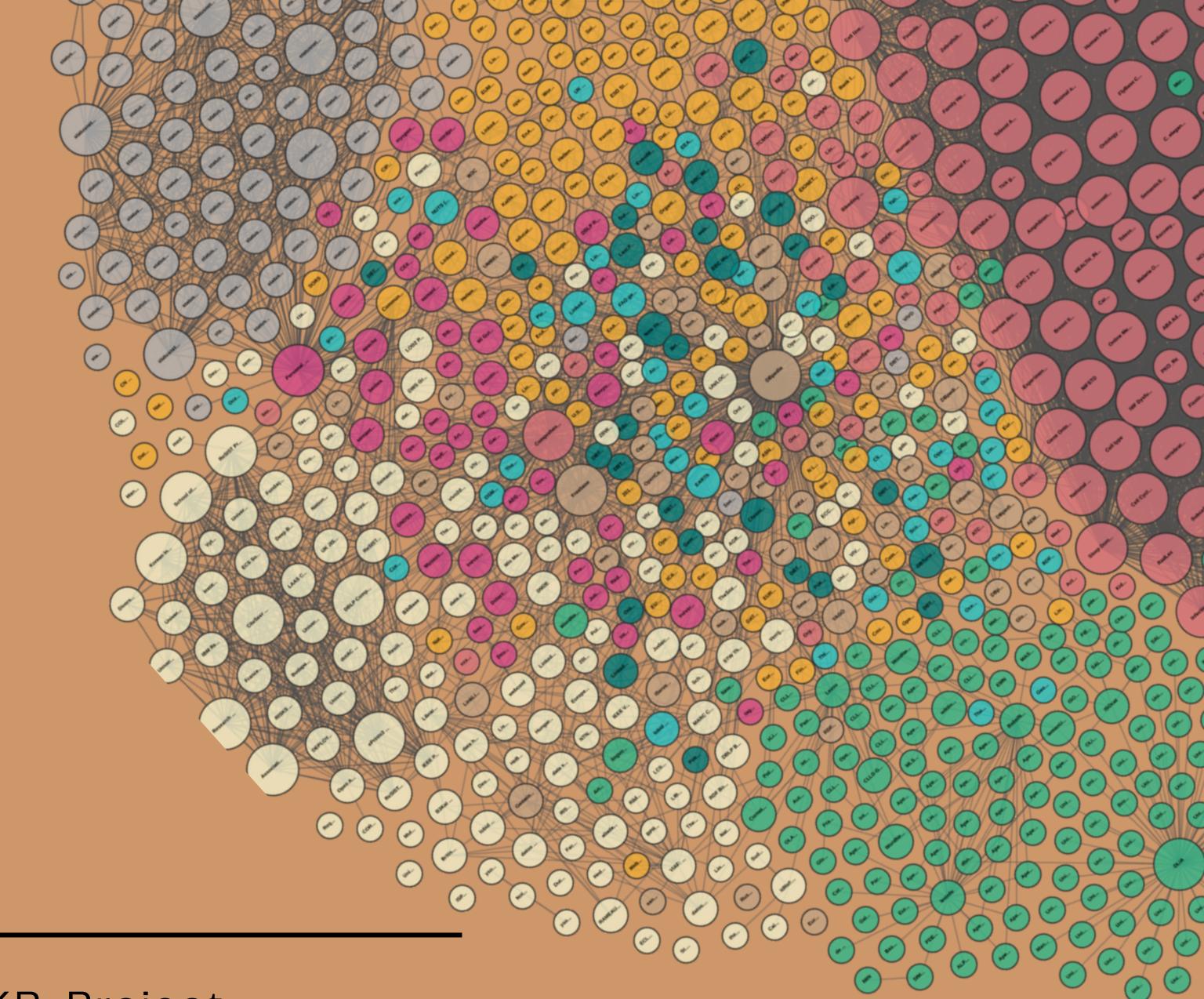
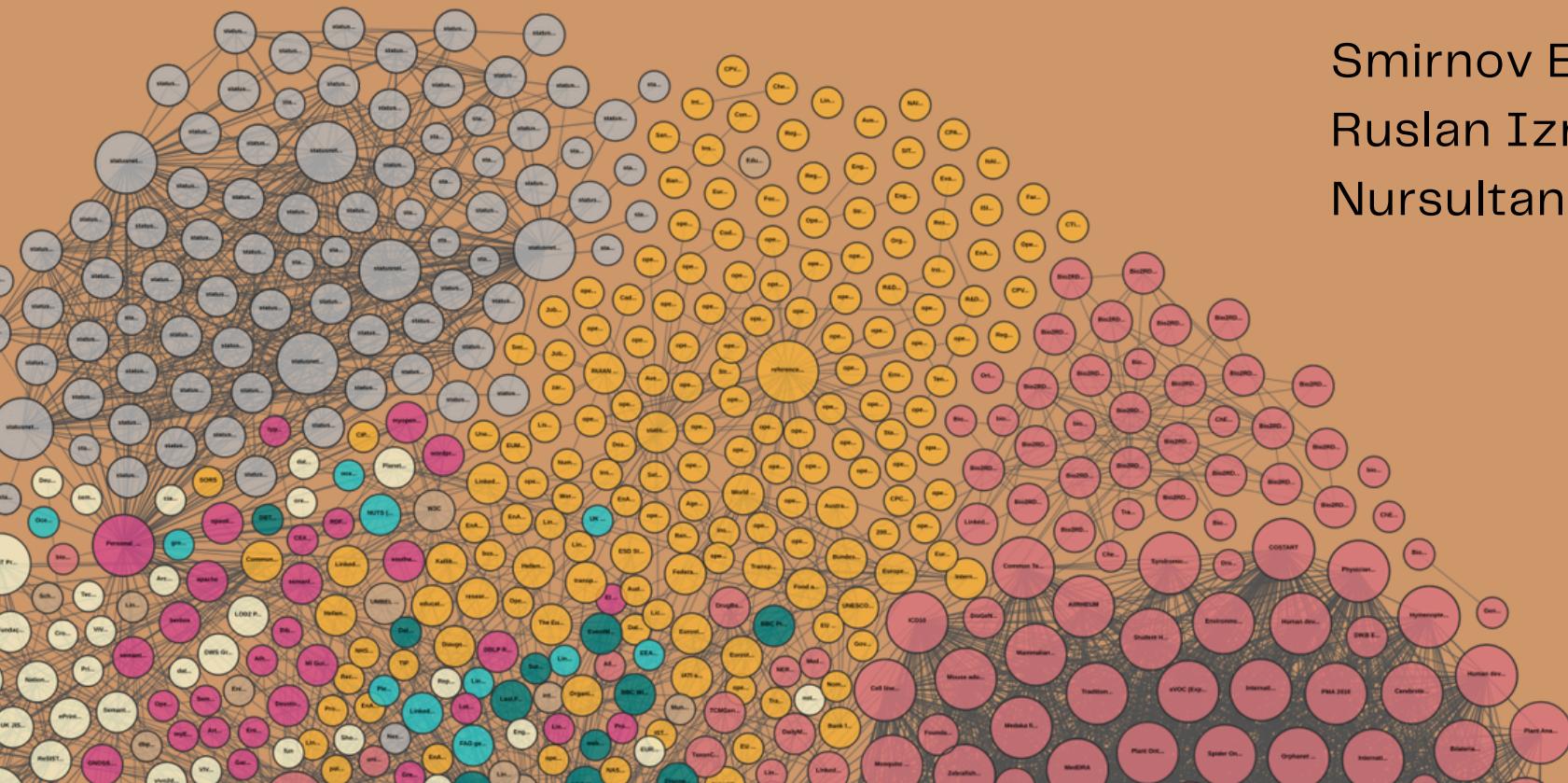
GRAPH TRANSFORMERS

github.com/Hexy00123/F24-DKR-Project

Smirnov Elisei

Ruslan Izmailov

Nursultan Abdullaev



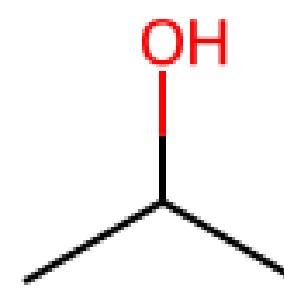
Problem Statement

Challenge:

Identifying toxic molecules in chemical datasets

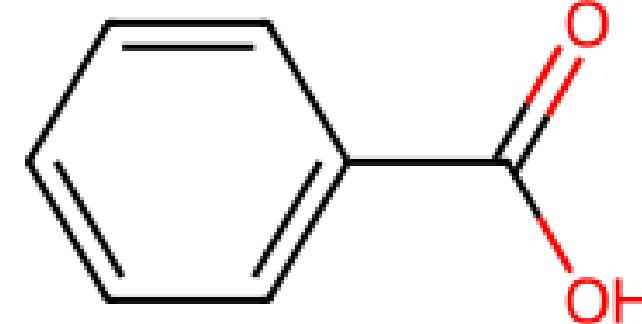
Why it Matters:

Essential for drug discovery and safety evaluation



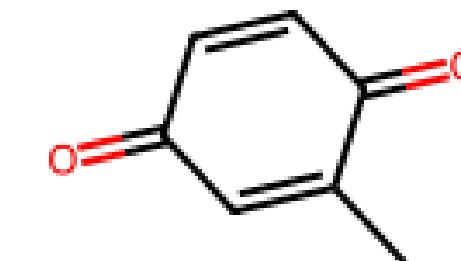
Non-Toxic

isopropanol



Non-Toxic

benzoic acid



Toxic

maleic anhydride

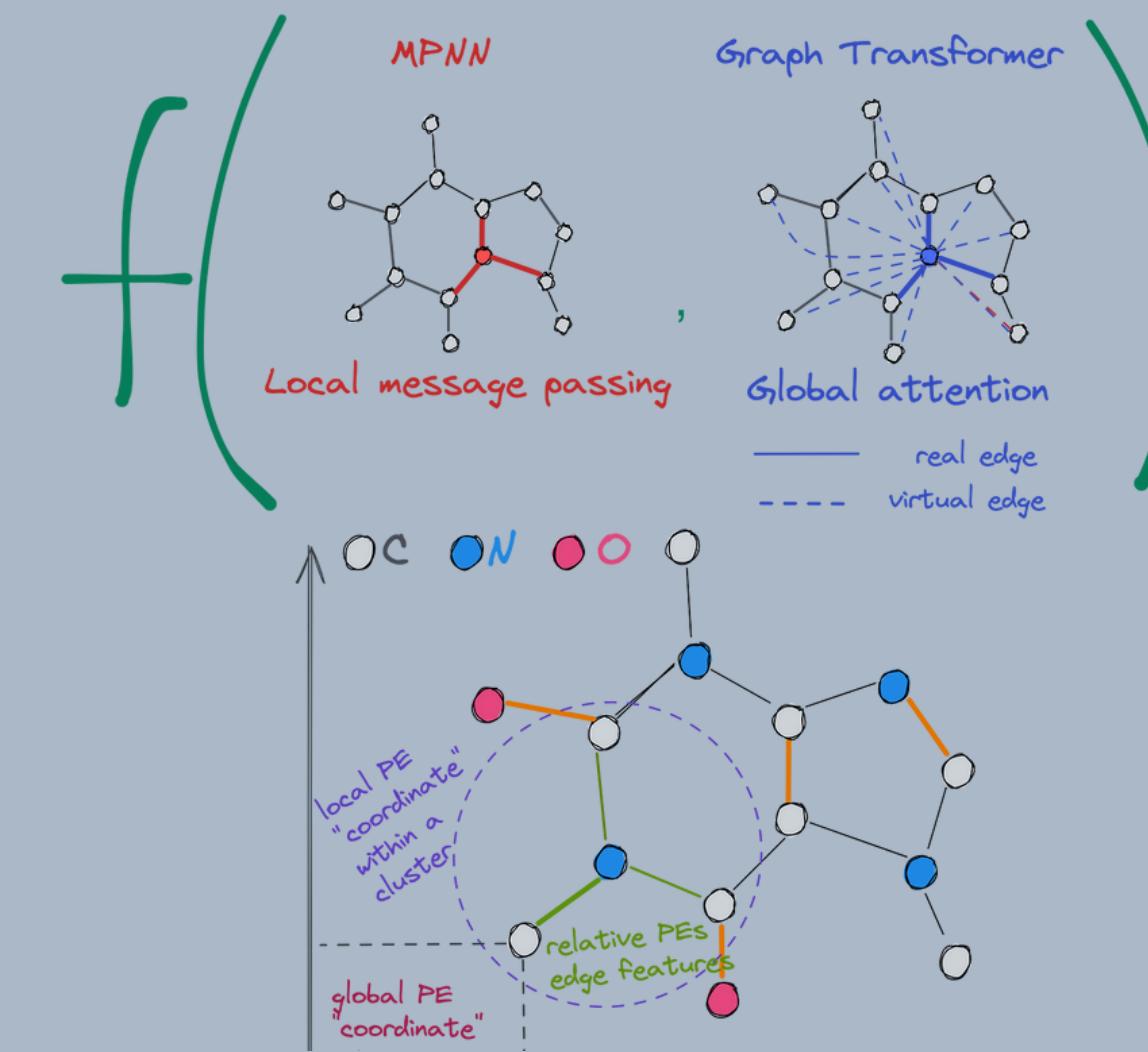
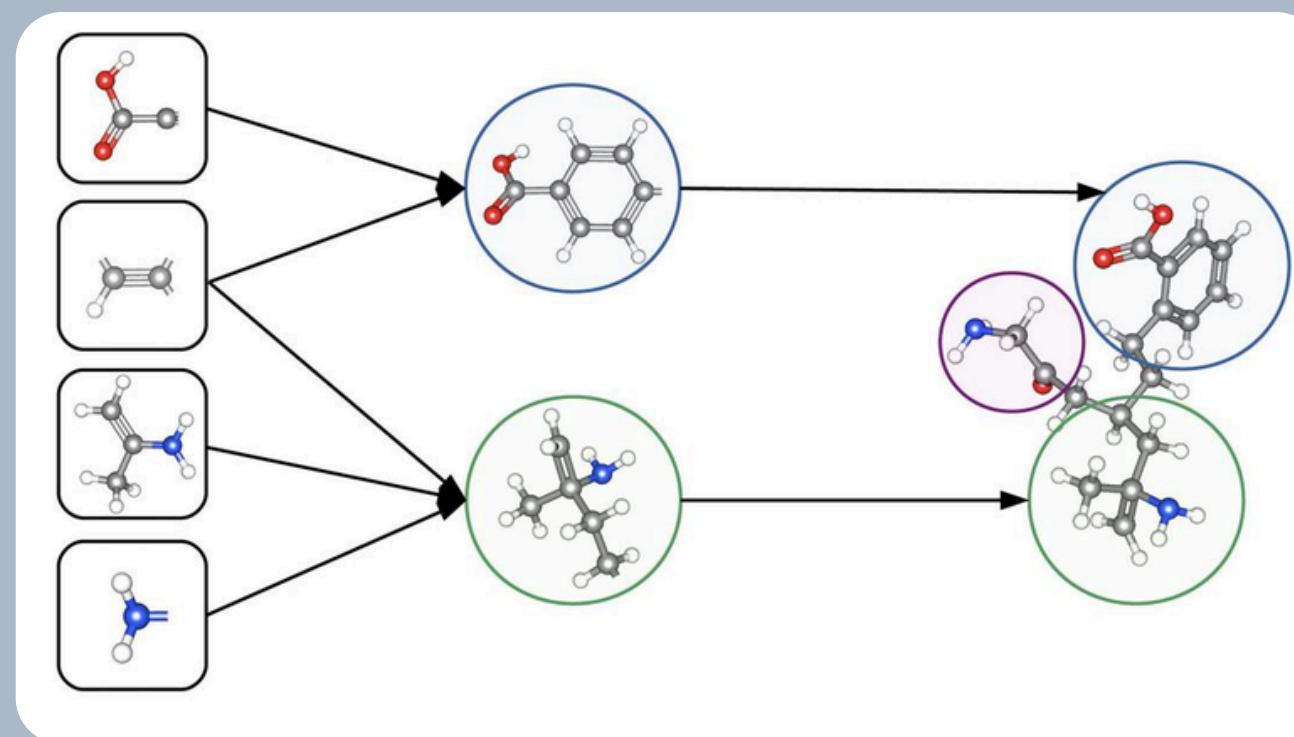


Toxic

selenourea derivative

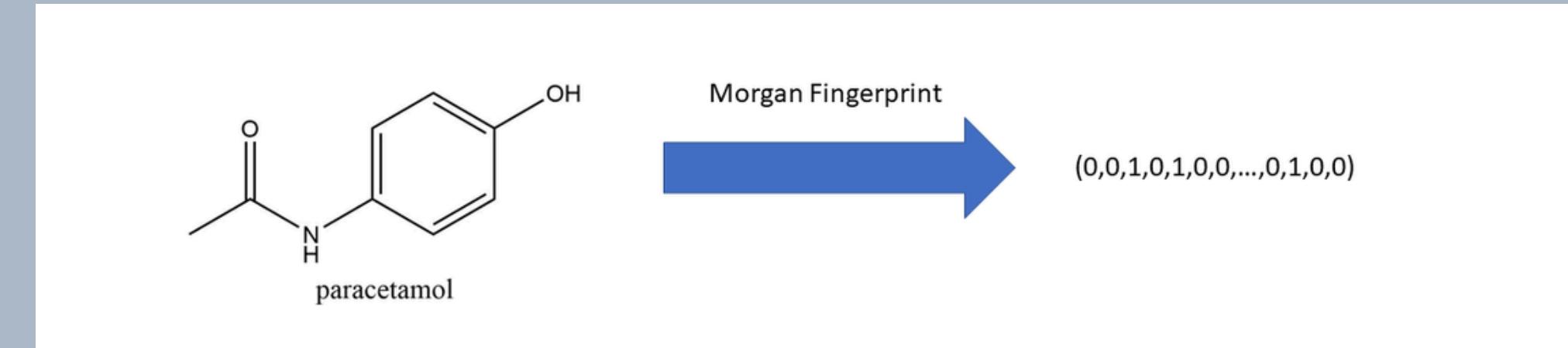
Methodology Overview

- 1) **Represent molecules as graphs** (atoms = nodes, bonds = edges) -> Tox21 dataset
- 2) **Extract features with Morgan fingerprints and adjacency matrices**
- 3) **Train a Graph Transformer model on toxicity prediction**



Extract features

1) **Morgan fingerprints** – are type of molecular descriptor used in computational chemistry.



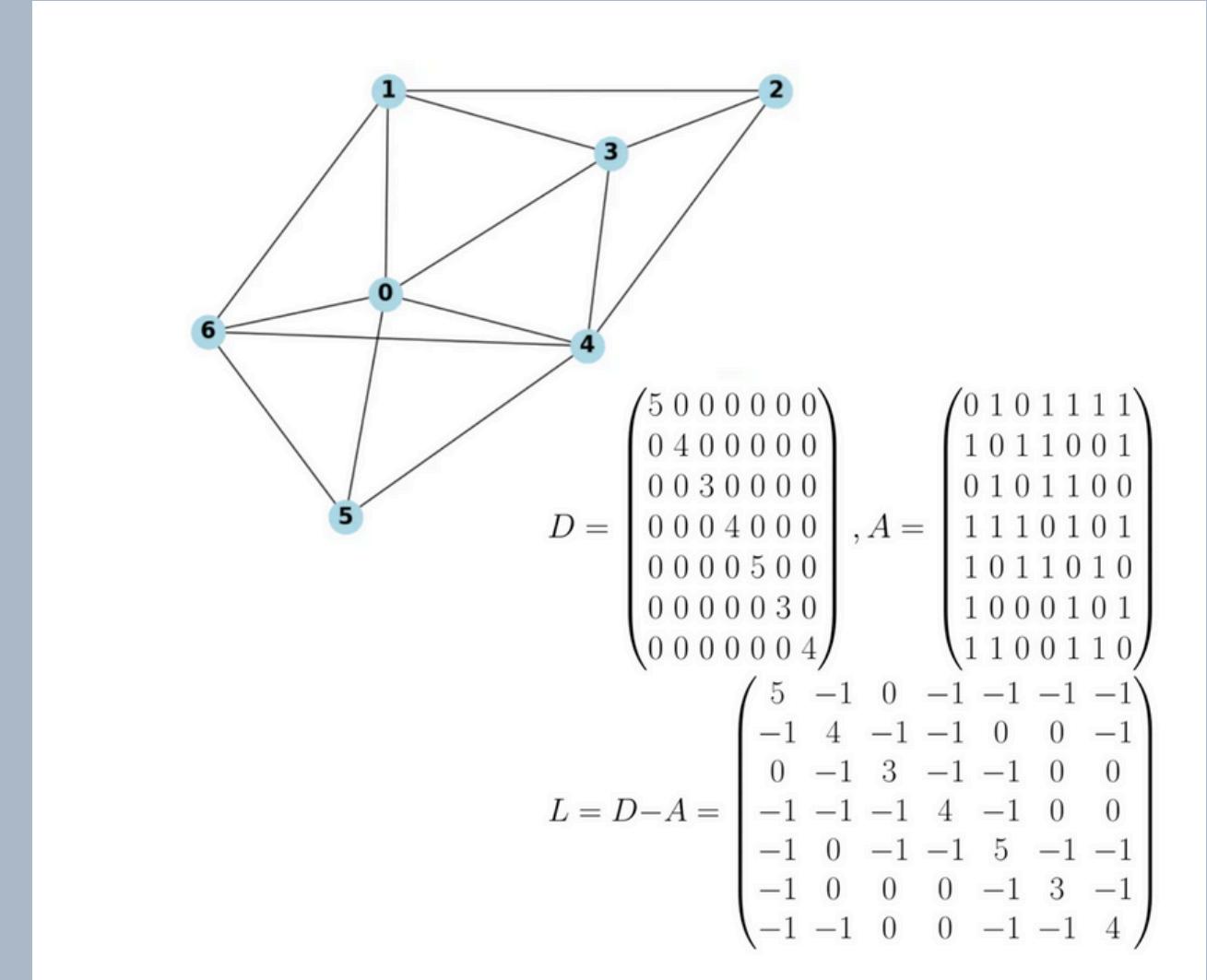
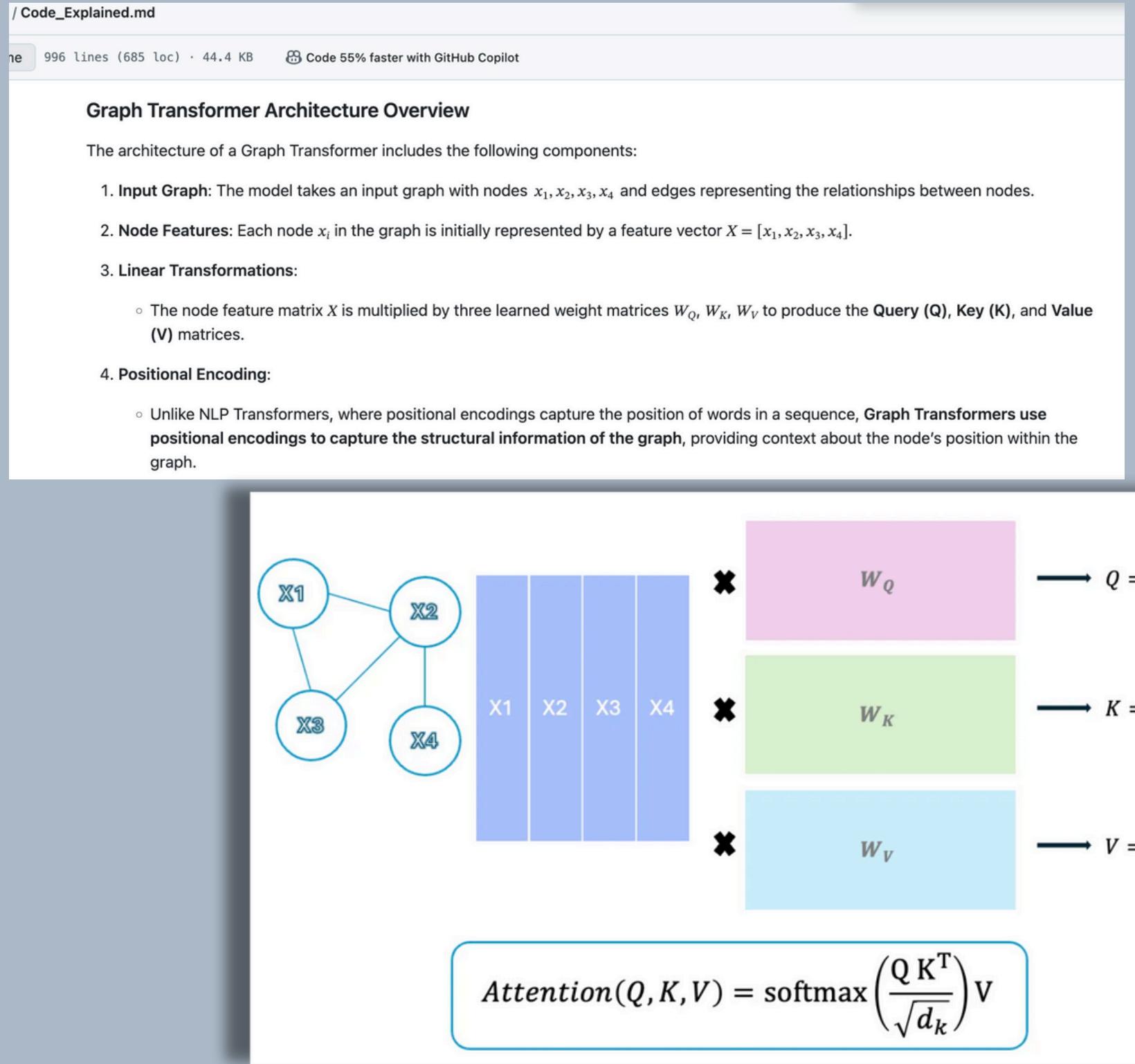
If the feature is present, the bit is set to 1; otherwise, it is 0.

2) **Adjacency Matrix**

Adjacency Matrix of Benzene						
0 -	0	1	0	0	0	1
1 -	1	0	1	0	0	0
2 -	0	1	0	1	0	0
3 -	0	0	1	0	1	0
4 -	0	0	0	1	0	1
5 -	1	0	0	0	1	0
	0	1	2	3	4	5

The matrix entry $\text{Adj}[i][j]$ is 1 if there is a bond between atom i and atom j, and 0 otherwise.

Graph Transformer model



Positional Encoding of Nodes Using
Laplacian Eigenvectors

Results

- Accurate toxicity prediction on Tox21 dataset is 70% with 2 epochs trained only!
- Graph Transformers excelled at capturing molecular interactions

```
Training: 100%|██████████| 735/735 [00:56<00:00, 13.08it/s]
Validating: 100%|██████████| 5873/5873 [00:20<00:00, 284.22it/s]
Validating: 100%|██████████| 1958/1958 [00:06<00:00, 289.37it/s]
Epoch 1/2: Loss: 0.6668, Train Acc: 0.6620, Test Acc: 0.6879
Training: 100%|██████████| 735/735 [00:57<00:00, 12.86it/s]
Validating: 100%|██████████| 5873/5873 [00:21<00:00, 273.15it/s]
Validating: 100%|██████████| 1958/1958 [00:06<00:00, 290.91it/s]
Epoch 2/2: Loss: 0.6346, Train Acc: 0.6874, Test Acc: 0.7074
```

```
new_smiles = "0" # water
prediction = predict(model, feature_extractor, new_smiles, device=device)
print(f"The predicted toxicity of the molecule with SMILES {new_smiles} is: {prediction}")
```

The predicted toxicity of the molecule with SMILES 0 is: 0