# Interpretable Graph Classification
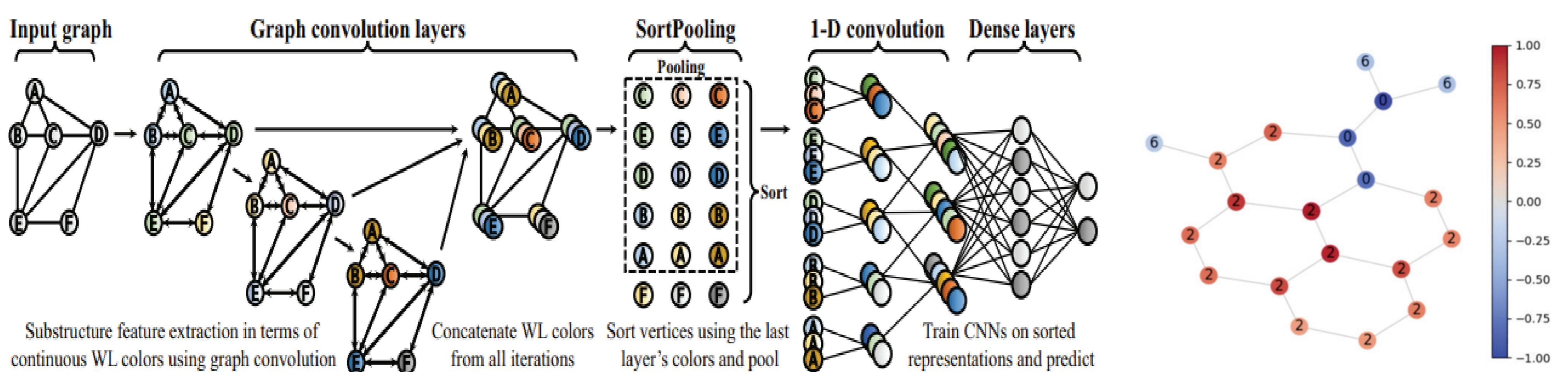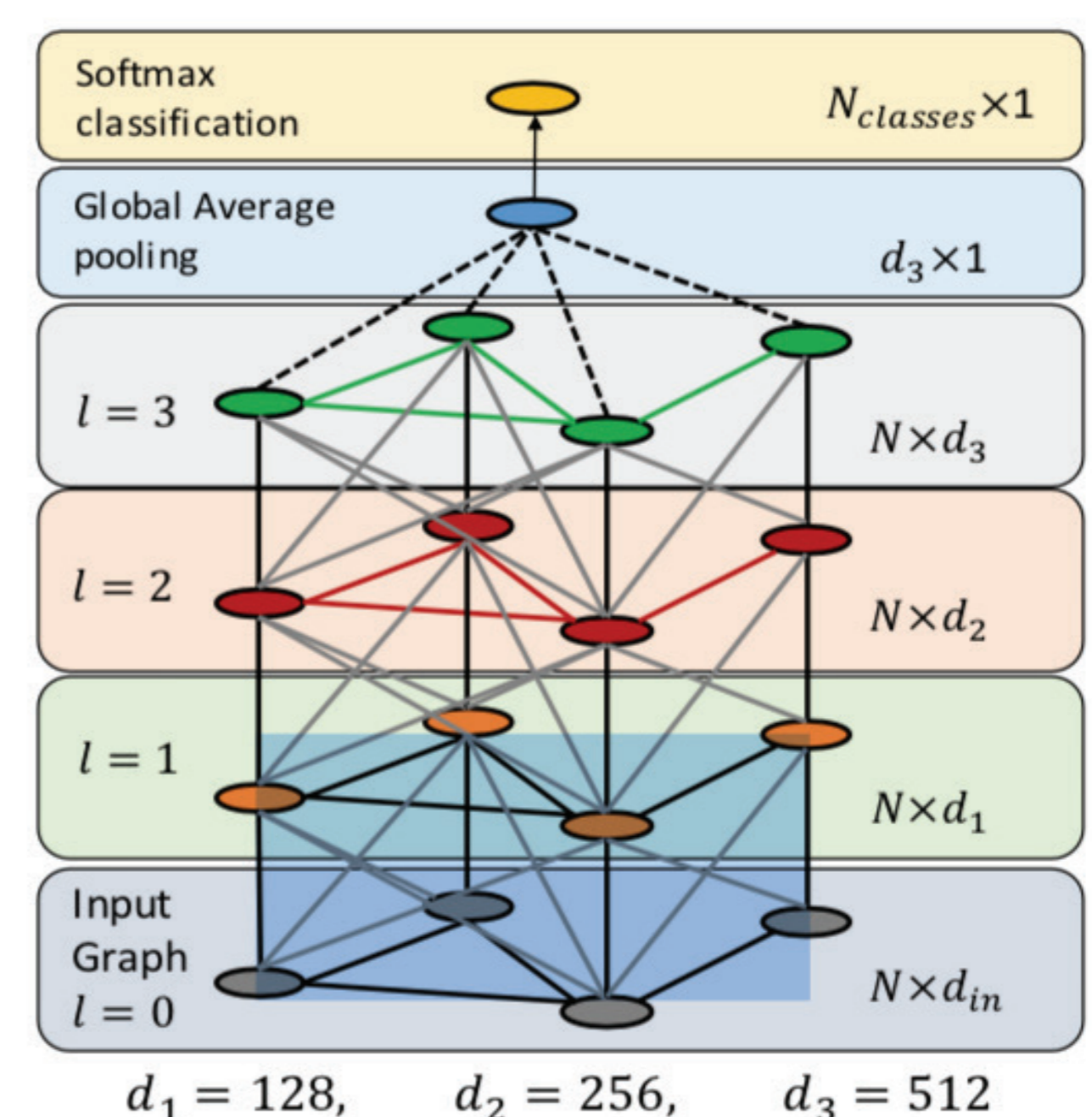
## Evaluation and interpretation of graph classification using deep learning

**Student:** Eko Edita Limanta          **Supervisor:** Asst Prof Arijit Khan

Neural networks are black-box in nature. In many real-life applications, it is of vital importance to know why the neural networks make certain decisions. For example, if a protein is classified as toxic due to presence of certain molecules.

This project aims to evaluate several interpretability methods for graph classification using deep learning. Two deep learning architectures are used: Deep Graph Convolutional Neural Network (DGCNN) and Graph Convolutional Networks(GCN).



$$d_1 = 128, \quad d_2 = 256, \quad d_3 = 512$$



Two interpretability methods – DeepLIFT and Contrastive Gradient – are then applied to these architectures. The interpretability results are then evaluated using three quantitative metrics: fidelity, contrastivity, and sparsity.