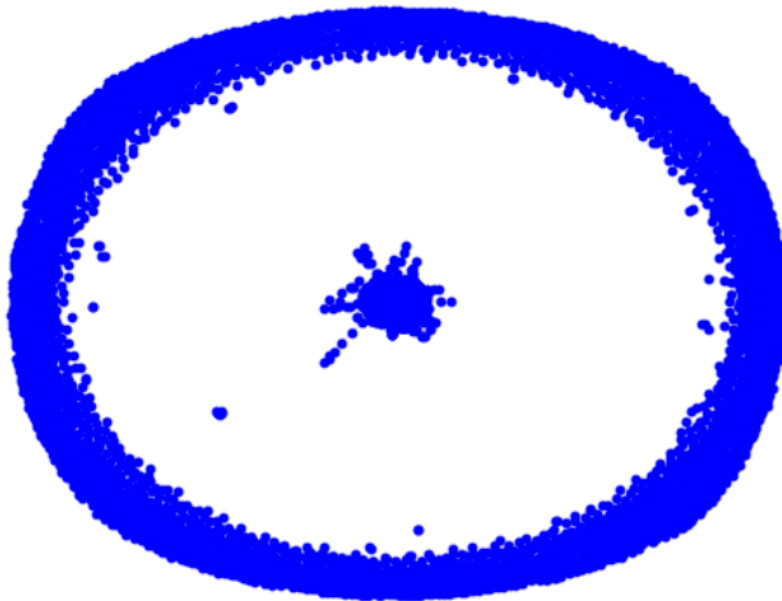**BRAUDE**
College of Engineering, Karmiel

Software Engineering Department
Braude College

# Graph Generation using Graph U-Nets

Capstone Project Phase B – Course 61999

24-1-R-4

Yana Raitsin 316086941
Yana.Raitsin@e.braude.ac.il

Hanil Zarbailov 319331419
Hanilz1995@gmail.com

Supervisor:
Dr. Renata Avros

# Table of Contents

**Abstract**

In the rapidly evolving landscape of graph representation learning and generative modeling, the quest for effective methodologies to understand and generate complex graph structures remains a pivotal challenge. This project introduces a model that integrates two architectures: MISC-GAN and Graph U-Nets. Our proposed model offers a unified framework for graph generation, leveraging the strengths of both models.

Index Terms: Generative adversarial network; Coarse graphs; Granularity levels; U-Nets; Cycle Consistency; Cycle-GAN; Multi-scale Analysis; Hierarchical Graph Representation.

# 1. Introduction

## 1.1 Problem statement

While the Misc-GAN model effectively captures the overall structure of input graphs, its "Multi-Scale Graph Representation Module" has limitations in precision. The proposal involves integrating a trained Graph U-Net into this module after the extraction of the coarse graphs. Graph U-Nets are adept at handling graph-structured data and can enhance the extraction of coarse graphs at various granularity levels. This integration aims to improve the accuracy of graph generation in the Misc-GAN model, leading to better performance in complex graph data analysis.

## 1.2 Motivation for Empowering Misc-GAN using Graph U-Nets

Graph U-Nets, inspired by the success of U-Nets in image processing, offer a promising direction for enhancing Misc-GAN. Graph U-Nets are adept at capturing hierarchical and topological features of graphs, essential for understanding complex network structures.

Misc-GAN is used to generate synthetic graphs in each granularity layer and then merge them all to construct a synthetic graph that is very similar to the source graph. The accuracy and relevance of these granularity layers heavily depend on the coarse graphs extracted from the original graph in the "Multi-Scale Graph Representation Module". Enhancing this module's outcome with Graph U-Nets is expected to greatly improve the accuracy of the synthetic graphs generated.

Graph U-Nets are particularly adept at considering node features in the source graph while preserving spatial information during graph reconstruction. This capability is pivotal for our enhancement strategy. By utilizing the strengths of Graph U-Nets to enhance and reconstruct coarse graphs, a significant improvement in the quality of synthetic graphs generated by Misc-GAN is anticipated.

## 1.3 Special uses

The improved Misc-GAN model, augmented with Graph U-Nets, opens up a plethora of applications that were previously challenging or inefficient with the standard Misc-GAN. This enhanced model is particularly valuable in domains where understanding and replicating complex graph structures are crucial.

One prominent area of application is in social network analysis. The enhanced model can effectively replicate and study social networks' dynamics, accounting for intricate community structures and evolving relationships. This is crucial in understanding social phenomena, network growth patterns, and information dissemination.

Another critical application is in biological networks, particularly in protein-protein interaction networks and neural network mapping. The model's ability to capture multi-scale structures aids in unraveling the complexities of biological systems, facilitating advancements in drug discovery and understanding of neural pathways.

Furthermore, the model finds relevance in technological networks, like the Internet and communication networks. It can simulate network behaviors under different scenarios, aiding in robust network design and analysis of network vulnerabilities.

In summary, the integration of Graph U-Nets into Misc-GAN expands its utility for analyzing and synthesizing complex networks across various domains. This enhanced model is not only a step forward in generative models for graphs but also a bridge connecting intricate network structures to practical, real-world applications.

## 2. Background

### 2.1 Neural Network

A neural network represents a paradigm within the field of artificial intelligence, wherein computational systems assimilate data through a mechanism inspired by the neural architecture of the human brain. Specifically categorized as a subset of machine learning known as deep learning, neural networks leverage an arrangement of interconnected nodes or neurons organized in layers, mirroring the organizational structure of the human brain. This configuration engenders an adaptive system, enabling computers to iteratively refine their performance by learning from errors. Consequently, artificial neural networks are applied to address intricate problem domains, exemplified by endeavors such as document summarization and facial recognition, with a heightened emphasis on precision.

The training process of a neural network entails the provision of input data and subsequently adjusting its weights to mitigate the disparity between the anticipated output and the observed actual output.

### 2.1.1 Backpropagation

Backpropagation, also recognized as the backward propagation of errors, constitutes a method for error evaluation that proceeds in reverse from output nodes to input nodes. It serves as a valuable mathematical instrument for enhancing predictive accuracy in the domains of data mining and machine learning. Neural networks employ backpropagation to calculate a gradient descent with respect to weight values across varied inputs. The systems undergo refinement through the adjustment of connection weights, aiming to minimize the discordance between desired and actual system outputs to the greatest extent possible.

### 2.1.2 Loss function

The loss function in a neural network calculates the disparity between the current output produced by the algorithm and the expected output. This metric is pivotal in evaluating the algorithm's efficacy in modeling the underlying data.

### 2.1.3 Activation function

An activation function plays a pivotal role in determining the activation state of a neuron within a neural network, serving as a decision mechanism to ascertain the significance of the neuron's input during the predictive process. By employing straightforward mathematical operations, the activation function derives the output based on a given set of input values supplied to a node or layer.

### 2.1.4 Bias in Neural Networks

Bias in neural networks is characterized as a constant term added to the product of features and weights, serving the purpose of offsetting the result. This incorporation of bias facilitates the models in adjusting the activation function, thereby influencing its shift towards either the positive or negative side.

### 2.1.5 Convolutional Neural Network

Convolutional Neural Networks (CNNs) represent a specialized class of neural networks designed for the processing of grid-like data structures, notably applied in contexts such as image analysis and recognition. Primarily employed for tasks involving visual information, CNNs excel in learning intricate visual features and patterns through the examination of pixel correlations, a process referred to as feature extraction. Their applicability extends to picture classification, object detection, and segmentation, where the hierarchical learning of features is achieved by the application of numerous filters to the input data. These filters enable the network to discern patterns and features at varying levels of abstraction, encompassing low-level attributes such as edges and corners, to more complex high-level features, including shapes and objects.

### 2.1.6 Convolutional Layer

The fundamental computational operation within a Convolutional Neural Network (CNN) involves a series of feature detectors, typically represented by a 3x3 matrix, traversing the input data. This process applies a mathematical operation known as convolution, wherein the feature detectors compute the dot product with localized regions of the input data. The outcome of this convolution is a collection of feature maps or activation maps, which effectively encapsulate distinct local patterns or features inherent in the input data. Convolutional layers play a pivotal role in extracting meaningful features across diverse scales and orientations, contributing to the network's capacity for hierarchical pattern recognition.

### 2.1.7 Pooling Layer

Pooling layers, recognized as a form of downsampling, orchestrate the reduction of dimensionality, thereby curtailing the number of parameters within the input. Analogous to the convolutional layer, the pooling operation involves the traversal of a filter across the entire input; however, the critical distinction lies in the absence of weights associated with this filter. Instead, the kernel applies an aggregation function to values within the receptive field, culminating in the generation of the output array. Although the pooling layer results in a loss of information, it engenders several advantages within Convolutional Neural Networks (CNNs), including complexity reduction, enhanced computational efficiency, and mitigation of overfitting risks.

### 2.1.8 Max Pooling

Max Pooling is a pooling operation that calculates the maximum value for patches of a feature map and uses it to create a downsampled (pooled) feature map. It is usually used after a convolutional layer. It adds a small amount of translation invariance - meaning translating the image by a small amount does not significantly affect the values of most pooled outputs.

### 2.1.9 Fully Connected Layer

The fully connected layer in a Convolutional Neural Network plays a pivotal role in global feature learning by connecting all neurons from the preceding layer to each neuron in the subsequent layer. Unlike convolutional and pooling layers that focus locally, the fully connected layer captures complex relationships across the entire input space. It forms the final stage of the network, learning weights for connections during training and consolidating information for comprehensive predictions or classifications. This layer's adaptability in learning intricate relationships makes it crucial for tasks requiring a holistic understanding of input data.

### 2.1.10 Graph Convolutional Network

Graph Convolutional Network is a specialized neural network architecture designed for processing graph-structured data, commonly found in applications such as social networks and molecular graphs. Unlike traditional Convolutional Neural Networks designed for regular grids like images, GCNs extend convolutional operations to irregular graph structures. They address the challenge posed by the non-grid topology of graphs by employing a message-passing approach, where each node aggregates information from its neighbors, simulating a convolutional operation. This enables GCNs to effectively capture and integrate information from the graph topology, making them suitable for tasks where understanding relational dependencies within the data is crucial.

### 2.1.11 Message Passing

Within the architecture of Graph U-Nets, message passing constitutes a fundamental mechanism for information propagation and aggregation across the graph-structured data. In the encoder phase, nodes engage in message passing to exchange information with their neighboring nodes, thereby facilitating the dissemination of local contextual knowledge. Subsequently, the received messages undergo aggregation, commonly accomplished through techniques such as summation or weighted aggregation, resulting in a refined representation of each node's neighborhood that encapsulates both local and global features. This aggregated information is then employed to update the features of individual nodes. The U-Net structure further integrates this process with downsampling in the encoder phase, akin to pooling layers, and upsampling in the decoder phase. The strategic implementation of skip connections ensures the preservation of high-resolution information during the downsampling and upsampling operations. The orchestration of message passing in Graph U-Nets, therefore, enables the model to effectively capture intricate dependencies within the graph, rendering it adept at tasks such as point cloud segmentation or molecular property prediction.

### 2.1.12 Softmax

The Softmax activation function plays a crucial role in neural networks, particularly in the context of classification tasks. It is commonly used in the output layer of neural networks when dealing *with* categorical target variables.

**Concept and Functionality of Softmax:**
The Softmax function is a generalization of the logistic function.
It converts a vector of raw scores (also known as logits) from the network into a probability distribution.
The output probabilities of Softmax sum up to 1, making it suitable for representing probabilities across multiple classes.
Mathematical Representation:

Given a vector $z$ of raw scores for each class, the Softmax function is defined as:

$$Softmax(z)_i = \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$$

where $i$ is the index of a particular class, $K$ is the total number of classes, $z_i$ is the raw score for class $i$, and $e$ is the base of the natural logarithm.

The numerator, $e^{z_i}$, is the exponentiated score for class $i$.

The denominator $\sum_{j=1}^{K} e^{z_j}$ is the sum of exponentiated scores for all classes, ensuring that the probabilities sum up to 1.

**Role in Classification:**

In classification problems, Softmax is used to interpret the outputs of a neural network as probabilities.

Each output value represents the probability that the given input belongs to a particular class.

This is particularly useful for tasks where an input can belong to only one of many possible classes (e.g., image classification, where an image might represent one of several different objects).

**Advantages:**

Softmax provides a probability-based output for multi-class classification problems, offering a measure of confidence or certainty in the classification.

It helps in determining the most likely class for a given input.

In summary, the Softmax activation function is essential for transforming the outputs of a neural network into a probability distribution over predicted output classes. This function is particularly beneficial in scenarios requiring a clear distinction between multiple categorical outcomes.


## 2.2 U-Nets

U-Nets are a type of neural network architecture mainly used in image processing, particularly for tasks like image segmentation. The main characteristic of U-Nets is their architecture, which is shaped like the letter "U". This structure consists of two main paths: an encoding (or contraction) path and a decoding (or expansion) path.


### 2.2.1 Encoding and Decoding in U-Nets

The concepts of encoding and decoding are fundamental in the architecture of U-Nets, playing a critical role in their functionality. Encoding in U-Nets refers to the process of compressing the input data into a compact representation, capturing the essential features while reducing dimensionality. This process involves a series of convolutional and pooling layers that progressively downsample the input, abstracting and retaining critical information while minimizing data volume.

Decoding, on the other hand, is the process of reconstructing the data from its encoded form. This step is crucial for generating output that is comparable in dimension and detail to the original input. In U-Nets, decoding involves a series of up-sampling and convolutional layers that gradually upscale the encoded representation back to its original size, while refining details and features. This symmetrical structure enables the network to efficiently learn and generate high-fidelity outputs, making U-Nets particularly effective for tasks requiring precise localization and detailed reconstruction.

### 2.2.2 Pooling in U-Nets

Pooling layers in U-Nets reduce the spatial dimensions of the input data. This reduction is crucial for minimizing computational complexity and for abstracting the input data into a more manageable form. The pooling operation, typically a max pooling, selects the most prominent features in a region, effectively condensing the information while maintaining the most critical aspects of the data. This process not only reduces the risk of overfitting by providing an abstracted form of the input but also aids in capturing invariant features, which are essential in many applications like image segmentation.

### 2.2.3 Unpooling in U-Nets

Unpooling layers in U-Nets are designed to reverse the effects of pooling, aiming to reconstruct the detailed structure from the condensed feature representation. This process involves up-sampling the feature maps to a higher resolution, which is crucial for tasks requiring detailed, high-resolution output. The unpooling operation often uses stored indices from the corresponding pooling layers to reconstruct the original structure of the data as accurately as possible. This approach ensures that the critical features lost during pooling are effectively recovered, enabling the U-Net to generate outputs that closely resemble the original input in terms of detail and structure.

### 2.3 Graph U-Nets

### 2.3.1 G-Pooling

Pooling is an operation which is used in CNN in order to get better generalization and performance and is best used on grid-like data such as images.
This operation takes the data and partitions it into non-overlapping rectangles, on which certain down-sampling functions are applied.
This operation cannot be applied to graphs, because there is no locality information among nodes in graphs, thus the partition operation is not applicable here. In general, the normal pooling operation outputs inconsistent results, thus, a new algorithm should take its place.
To improve the Pooling operation, the algorithm proposes the graph pooling (gPool) layer to enable down-sampling on graph data as follows:

| Inputs | Projection | Top k Node Selection | Gate | Outputs |

$$y = \frac{X^l p^l}{||p^l||}$$
$$idx = rank(y, k)$$
$$\tilde{y} = sigmoid(y, k)$$
$$X^{\sim l} = X^l(idx, :)$$
$$A^{l+1} = A^l(idx, idx)$$
$$X^{l+1} = \underline{X}^l \odot (\tilde{y} 1_C^T)$$

1. Projection and Scoring ($y = \frac{X^l p^l}{||p^l||}$)
   - $X^l$ is the feature matrix of the nodes in the graph.
   - $p^l$ is a trainable projection vector.
   - The projection $X^l p^l$ computes the scalar projection of each node's features onto $p$.
   - The division by $||p^l||$ (the norm of $p^l$) normalizes this projection.
   - The result $y$ is a vector of scores estimating the importance of each node based on its projection value.
2. Node Ranking and Selection ($idx = rank(y, k)$):
   - $rank(y, k)$ operation ranks the nodes based on their scores in $y$.
   - It selects the indices ($idx$) of the top $k$ nodes with the highest scores.
   - $k$ is a predefined number representing how many nodes to select.
3. Gating Mechanism ($\tilde{y} = sigmoid(y, k)$):
   - $\tilde{y}$ is the gated score vector obtained by applying the sigmoid function to the selected scores $y(idx)$.
   - Sigmoid function is used to scale these scores between 0 and 1.
4. Feature Matrix for Selected Nodes ($X^{\sim l} = X^l(idx, :)$):
   - $X^{\sim l}$ is the new feature matrix formed by extracting the rows from $X^l$ corresponding to the selected node indices $idx$.
5. New Adjacency Matrix ($A^{l+1} = A^l(idx, idx)$):
   - $A^{l+1}$ is the new adjacency matrix created by selecting the rows and columns from the original adjacency matrix $A^l$ corresponding to the selected nodes.

12

6. Node Feature Update ($X^{l+1} = \underline{X}^l \odot (\tilde{y}1_C^T)$):
   - $X^{l+1}$ is the updated feature matrix for the next layer.
   - $\odot$ denotes element-wise multiplication.
   - $\tilde{y}1_C^T$ is a vector used in the gating operation, where $1_C$ is a vector of size $C$ (number of features) with all components being 1.
   - This step controls the information flow from the selected nodes, emphasizing the features of more important nodes.

### 2.3.2 G-Unpooling

The gUnpool layer uses the position information extracted from the gPool operation before it to reconstruct the original graph structure by using empty feature vectors for unselected nodes.

The purpose of the operation performs the inverse operation of the gPool layer and restores the graph into its original structure, with empty feature vectors for unselected nodes.

$$X^{l+1} = distribute(0_{NxC}, X^l, idx)$$

It records the locations of nodes selected in the corresponding gPool layer and uses this information to place nodes back to their original positions in the graph.

$idx \in z^{*k}$ contains the indices of the selected nodes in the corresponding gPool layer that reduces the graph size from N nodes to k nodes.

$X^l \in R^{kxC}$ are the feature matrix of the current graph.

$0_{NxC}$ is the initially empty feature matrix for the new graph.

$distribute(0_{NxC}, X^l, idx)$ is the operation that distributes row vectors in $X^l$ into $0_{NxC}$ feature matrix according to their corresponding indices stored in $idx$.

In $X^{l+1}$ row vectors with indices in idx are updated by row vectors in $X^l$, while other row vectors remain zero.

### 2.4 GAN

Generative Adversarial Networks (GANs) consist of two neural networks, the **generator** and the **discriminator**, which are trained simultaneously. The generator creates data that is as close as possible to real data, while the discriminator evaluates the authenticity of both real and generated data. The networks engage in a min-max game where the generator tries to fool the discriminator, and the discriminator tries to accurately distinguish between real and generated data.

### 2.4.1 Generator

The generator serves as the creative component, tasked with generating new data instances that mimic the real data. It aims to produce outputs indistinguishable from actual data, effectively fooling the discriminator. The generator's effectiveness is honed over time through its adversarial relationship with the discriminator, learning to create increasingly convincing fakes.

### 2.4.2 Discriminator

The discriminator acts as the critic, differentiating between real and generated data. Its primary function is to evaluate the authenticity of the input, whether it is a genuine data instance or a fabrication from the generator. The discriminator's role is crucial in training the generator to produce more realistic data, as it provides feedback that guides the generator's improvements.

The GAN Formula: The objective function of GANs is given by:

$$min_G max_D V(D, G) = E_{x \sim P_{data}(x)}[log D(x)] + E_{z \sim P_z(z)}[log(1 - D(G(z)))]$$

In this formula, $D(x)$ is the discriminator's estimate of the probability that real data instance $x$ is real. $G(z)$ is the generator's output when given noise $z$. $P_{data}$ is the data's distribution and $P_z$ is the noise's distribution.

### 2.5 Misc-GAN

Misc-GAN (Multi-scale Generative Adversarial Network for Graphs) is a sophisticated neural network model designed for graph data. It stands out for its ability to capture and model the distribution of graph structures at multiple levels of granularity. This multi-scale approach is essential for understanding complex graph organizations, especially in real-world networks that exhibit hierarchical distribution over graph communities.

### 2.5.1 Granularity Levels

One of the critical aspects of Misc-GAN is its focus on granularity levels. The model recognizes that real-world networks are not flat but have intricate hierarchical structures. By modeling these structures at different granularity levels, Misc-GAN can capture both local and global patterns in the graph, which is vital for a more nuanced and accurate representation of the graph data.

### 2.5.2 Architecture of Misc-GAN

The architecture of Misc-GAN is built upon three main modules, each contributing to its unique ability to generate and reconstruct graph data:

1) **Multi-Scale Graph Representation Module:** This module is the first stage of Misc-GAN, where the target domain graph is taken as input. The module explores the hierarchical structures within the graph to extract important structures at each granularity layer. It creates a series of coarse graphs, each representing the graph at different levels of granularity.

2) **Graph Generation Module:** At this stage, Misc-GAN works at each granularity level to create synthetic coarse graphs. This process involves dividing the original graph into subgraphs using advanced clustering techniques and then generating block diagonal matrices. These matrices capture community-level graph structures, allowing the model to understand and represent the graph at different levels effectively.

3) **Graph Reconstruction Module:** The final stage involves combining all synthetic coarse graphs created in the previous stage to generate the final graph. This process uses a linear reconstruction function that preserves the hierarchical structures at different granularity levels, ensuring the scalability and efficiency of the model.

### 2.5.3 Cycle Consistency

Cycle consistency is integral to the Misc-GAN architecture. It ensures the graph generation process is reliable and the reconstructed graph maintains the original graph's essential properties. This concept is borrowed from Cycle-Consistent Adversarial Networks, where a generated graph is expected to be transformed back to its original form, thus ensuring that the generative process preserves critical graph properties.

### 2.5.4 Multi-Scale Analysis

The multi-scale analysis in Misc-GAN allows the model to capture and process graph data at multiple resolutions. This approach is beneficial for handling large-scale graphs that have numerous representations, making the model versatile and applicable to various types of graph data. The multi-scale approach also aids in reducing computational complexity while preserving important topological features of the graph.

### 2.5.5 Cycle Gan

**(1)** This concept relies on GAN but duplicated to achieve more consistent results and better mapping from each domain to the other.

**(2)** We use 2 discriminators and 2 generators: $D_x$, $D_y$, $G: X \to Y$ and $F: X \to Y$ (while $X$ and $Y$ are two domains that depict the source graph and the generated graph domains).

**(3)** Each set of discriminators and generators transfer from one domain to the other, very similarly to GAN.

To view this concept in action, refer to the next section.

### 2.5.6 A Generic Joint Learning Framework

The framework is presented as a generic optimization problem with a multi-objective function. This function comprises four key components: multi-scale reconstruction loss, forward adversarial loss, backward adversarial loss, and cycle consistency loss (cycleGAN).

$$L = L_{ms} + L_F + L_B + L_{cyc}$$

$$= KL(\sum_{l=1}^{L} w^{(l)} F^{(l)}(G_s^{(l)}) + b, G_t)$$

$$+ \alpha \sum_{l=1}^{L} E_{G_t^{(l)} - P_{data}(G_t^{(l)})}[log D_F^{(l)}(G_t^{(l)})] + E_{G_s^{(l)} - P_{data}((G_s^{(l)}))}[log(1$$

$$- D_F^{(l)}(F^{(l)}(G_s^{(l)})))]$$

$$+ \beta \sum_{l=1}^{L} E_{G_s^{(l)} - P_{data}(G_s^{(l)})}[log D_B^{(l)}(G_s^{(l)})] + E_{G_t^{(l)} - P_{data}(G_t^{(l)})}[log(1 - D_B^{(l)}(B^{(l)}(G_t^{(l)})))]$$

$$+ \gamma \sum_{l=1}^{L} E_{G_s^{(l)} - P_{data}(G_s^{(l)})}[||B^{(l)}(F^{(l)}(G_s^{(l)})) - G_s^{(l)}||_1]$$

$$+ E_{G_t^{(l)} - P_{data}(G_t^{(l)})}[||F^{(l)}(B^{(l)}(G_t^{(l)})) - G_t^{(l)}||_1]$$

**Multi-Scale Reconstruction Loss ($L_{ms}$):** This term is designed to minimize the divergence between the target graph ($G_t$) and the generated graph ($G_{\tilde{t}}$). It uses a Kullback-Leibler ($KL$) divergence measure to compare two graphs. The divergence is calculated between the adjacency matrices of $G_t$ and $G_{\tilde{t}}$.

**Forward Adversarial Loss ($L_F$):** This component learns a forward mapping function from the source graph ($G_s$) to the target graph ($G_t$). It involves a discriminator ($D_F^{(l)}$) that distinguishes whether a graph is from the target domain or generated by the model.

**Backward Adversarial Loss ($L_B$):** Similar to $L_F$, but in reverse. It defines an adversarial loss that learns the mapping from the target domain to the source domain, with its discriminator ($D_B^{(l)}$).

**Cycle Consistency Loss** ($L_{cyc}$)**:** This term ensures that the mapping functions do not contradict each other by maintaining consistency in both directions of the mapping.

**Overall Objective:** The overall objective ($L$) is a sum of these four terms, with positive constants ($\alpha$, $\beta$, $\gamma$) balancing their impact. The model aims to optimize this objective function using a min-max game typical of GAN frameworks.

### 3. Solution – Graph Generation using Graph U-Nets.

The project aims to simulate graphs that closely resemble the structure and relationships of a given citation dataset by training a Graph U-Nets network on a related dataset. Specifically, we aim to replicate citation datasets such as CORA and PubMed. This network is later utilized for extracting coarse graphs and reconstructing the generated graph at each granularity layer into the final output graph.

The primary workflow of the model is as follows:

1. **Pre-training**: A Graph U-Nets network is pretrained on a dataset that is similar in nature to a citation dataset (such as DBLP and PROTEINS).
2. **Data Preprocessing**: The target dataset (such as CORA and PubMed) is processed to extract essential data and prepare it for the model.
3. **Coarse Graph Extraction**: Coarse graphs for each granularity layer are derived using the G-Pool operation of the trained Graph U-Nets.
4. **Coarse Graph Generation**: Cycle-GAN is applied to generate coarse graphs that mimic the original coarse graphs for each respective layer.
5. **Graph Reconstruction**: Two methodologies are evaluated:
   1. The trained network applies the G-Unpool operation on each synthetic coarse graph to replicate the structure of the original graph more accurately using node features. All final coarse graphs are then integrated into a final synthetic graph using the Graph Reconstruction Module.
   2. The Graph Reconstruction Module first merges all synthetic coarse graphs to form a comprehensive synthetic graph. The trained network's G-Unpool operation is then executed on this synthetic graph to reconstruct the structure of the original graph more faithfully using node features.

**3.1 Data Preprocessing**

In the integration of Misc-GAN and Graph U-Nets, data preprocessing is a crucial step that ensures raw data is transformed into a format that is suitable for model training and testing. This section outlines the preprocessing techniques used for both Misc-GAN and Graph U-Nets.

**3.1.1 Misc-GAN**

For Misc-GAN, the preprocessing process involves several key steps to prepare the dataset. Here's a detailed breakdown:

1. **Loading and Initial Setup**: The dataset is loaded from text files, and Google Drive is mounted to access the required data files.
2. **Supernode Creation**: The supernode function aggregates nodes into supernodes based on a specified clustering. This function iterates through clusters, summing the connections within each cluster to form supernodes, thus simplifying the graph and reducing its complexity.
3. **Removing Zero Rows**: The remove_zero_row function eliminates any rows and columns that sum to zero, which can result from isolated nodes or other preprocessing steps.
4. **AMG Clustering**: The AMG function performs algebraic multigrid (AMG) clustering, recursively coarsening the graph to create a hierarchy of coarser graphs. The fine2coarse function is central to this, selecting coarse nodes and computing interpolation matrices.
5. **DBSCAN Clustering**: The DBSCAN function applies the DBSCAN algorithm to cluster the data points. This density-based clustering algorithm identifies core samples, noise, and clusters.
6. **Edge Extraction**: The supernode function is used again to extract edges between supernodes, creating a simplified representation of the graph.
7. **Saving Preprocessed Data**: The processed data, including the adjacency matrix, interpolation matrices, and clustering results, is saved to .mat files for use with the Misc-GAN model.

These steps transform raw graph data into a structured format that Misc-GAN can utilize, reducing complexity while preserving essential graph properties.

### 3.1.2 Graph U-Nets

For Graph U-Nets, the preprocessing involves converting the dataset into a specific text format that the model requires. This process includes detailed information about the graph structure, nodes, and edges. The key steps are:

1. **Loading the Dataset**: The Google Drive is mounted to access the dataset files, which typically include graph indices, graph labels, and edge lists.
2. **Graph Index and Label Extraction**: The get_indic function reads the graph index file and extracts the graph IDs and their sizes. The get_labels function reads the label file and assigns labels to each graph.
3. **Graph Transformation**: The trans_graphs function writes the graph information to a text file in the required format. It reads the edge list file, organizes the edges by graph, and writes the graph size, label, and edge information.
4. **Detailed Node Information**: Within each graph block, nodes are described by their tags and neighbors. The write_graph function handles this, ensuring each node's neighbors are correctly indexed and listed.
5. **Loading Processed Data**: The load_data function reads the processed graph data file, constructs graph objects using NetworkX, and extracts node features and labels for use in Graph U-Nets.

The preprocessing steps for Graph U-Nets convert the raw dataset into a detailed text format that includes all necessary graph, node, and edge information. This structured format allows Graph U-Nets to effectively utilize the dataset for graph classification tasks.

### 3.2 Graph Generation using Graph U-Nets Algorithm

### 3.2.1 Method 1: G-Unpool after Generating Coarse Graph for Each Layer

In the first method of integrating Misc-GAN and Graph U-Nets, the G-Unpool operation is applied after generating the coarse graph at each layer. The process begins with the initial coarse graph generation, where the GPool operation is employed to coarsen the graph. This involves reducing the graph's complexity while retaining its essential structure and features.

After the coarse graph is generated for a layer, the G-Unpool operation is applied immediately. This operation reverses the pooling process to some extent by reconstructing the graph to a finer resolution. It leverages the previously stored downsampled features and indices to ensure that the reconstructed graph maintains the integrity and characteristics of the original graph.

The key steps in this method are as follows:

1. **Initial Coarse Graph Generation:** Use the GPool operation to coarsen the graph, reducing its size and complexity.
2. **G-Unpool Application:** Immediately after the coarse graph generation for each layer, apply the G-Unpool operation to partially reconstruct the graph. This involves using stored features and indices from the downsampling process.
3. **Iterative Process:** Repeat the above steps for each layer, ensuring that after each coarse graph generation, the G-Unpool operation is applied to maintain graph integrity.

This method ensures that the graph retains a higher level of detail throughout the process by reconstructing the graph after each layer's coarsening. It allows for better preservation of the graph's structure and features at each stage of the Misc-GAN process.


### 3.2.2 Method 2: G-Unpool after Reconstructing the Graph and Applying Sigma Operation

The second method involves applying the G-Unpool operation after the entire graph reconstruction and the sigma operation. This method follows a different sequence, where the graph is first coarsened layer by layer without immediate unpooling. Once the graph has been fully coarsened through all layers, it undergoes a reconstruction process using the sigma operation.

The key steps in this method are as follows:
1. **Layer-by-Layer Coarsening:** Use the GPool operation to coarsen the graph at each layer without immediately applying the G-Unpool operation.
2. **Graph Reconstruction:** After all layers have been processed, reconstruct the graph using the sigma operation. This step involves enhancing the graph's structure by applying a function that consolidates the features and connections within the graph.
3. **G-Unpool Application:** Once the graph reconstruction and sigma operation are complete, apply the G-Unpool operation to the fully reconstructed graph. This step ensures that the graph regains its finer resolution and detail, using the stored downsampled features and indices from the initial coarsening process.

This method focuses on maintaining the coarsened structure throughout the entire process, applying the unpooling only after the graph has been fully reconstructed. This approach can be beneficial for applications where the final reconstructed graph needs to be of high resolution and detail, but intermediate stages can afford to work with coarser representations.

### 3.3 Pre-Training

### 3.3.1 Training Graph U-Nets PROTEINS

Pre-training Graph U-Nets using the PROTEINS dataset is a crucial step for building a robust model capable of handling complex graph structures. The PROTEINS dataset is chosen for its suitability in graph classification tasks and to ensure the model can generalize effectively. This dataset helps prepare the model for handling diverse and intricate graph structures found in citation networks. The pre-training process involves data preparation, model initialization, training, evaluation, and model saving. Training on the PROTEINS dataset enables the model to develop a strong foundation for capturing the nuances of graph data, which is essential for citation network tasks.

### 3.3.2 Training Graph U-Nets DBLP

Pre-training Graph U-Nets with the DBLP dataset is essential for adapting the model to handle citation networks effectively. The DBLP dataset, representing academic citation graphs, is ideal for this purpose. This pre-training phase follows similar steps to those used with the PROTEINS dataset, including data preparation, model initialization, training, evaluation and model saving. By using the DBLP dataset, the model is well-prepared for applications in academic and research environments. The DBLP dataset's focus on citation networks ensures the model can accurately capture the relationships and patterns within such networks.

### 3.3.3 Training Graph U-Nets with COLLAB

The COLLAB dataset, which represents scientific collaboration networks, is used to further enhance the model's ability to understand complex graph structures. This pre-training phase also involves data preparation, model initialization, training, evaluation and model saving. Training on the COLLAB dataset ensures the Graph U-Nets model is robust and versatile, capable of handling various graph-structured data, including citation networks. The COLLAB dataset's emphasis on collaborative relationships makes it a valuable addition to the pre-training process, strengthening the model's ability to process and analyze citation networks effectively.

### 3.4 Simulation process with Graph Generation using Graph U-Nets

The simulation process using Misc-GAN aims to create synthetic citation networks that mimic real-world datasets. By modifying the algorithm and the original Misc-GAN approach, the realism and utility of the generated graphs are enhanced. This integration provides a robust framework for producing high-fidelity synthetic data, which can be utilized for various research and development purposes in citation networks. The following sections describe the simulation process with the CORA and PubMed datasets, including an evaluation of two G-Unpool methods.

### 3.4.1 Simulation with CORA

The CORA dataset, a well-known benchmark in citation network research, is used to validate the integration of Misc-GAN with Graph U-Nets.

**Data Preparation**: The CORA dataset undergoes preprocessing to extract relevant features, adjacency matrices, and labels, ensuring compatibility for training the Graph U-Nets model and generating synthetic graphs with Misc-GAN.

**Model Training**: The pre-trained Graph U-Nets model, fine-tuned on datasets such as PROTEINS, DBLP, and COLLAB, is further trained using the CORA dataset. This step enhances the model's understanding of the specific characteristics of citation networks in CORA.

**Synthetic Graph Generation**: Misc-GAN, integrated with the trained Graph U-Nets model, generates synthetic citation networks. G-Unpool operations are applied to reconstruct graph structures after coarsening steps, preserving the structural integrity and properties of real citation networks. Both G-Unpool methods are tested to evaluate their effectiveness in reconstructing the graphs.

**Evaluation**: The synthetic graphs are compared to the original CORA dataset using various graph metrics, such as node degree distribution and clustering coefficient, to ensure realism and utility for further research.

The CORA dataset simulation demonstrates the effectiveness of our approach in generating realistic citation networks. Insights gained from this process are applicable to other citation networks, validating the versatility and robustness of the Misc-GAN and Graph U-Nets integration.

### 3.4.2 Simulation with PubMed

The PubMed dataset, another key benchmark in citation network research, is used to further assess the capabilities of the Misc-GAN and Graph U-Nets integration.

**Data Preparation**: The PubMed dataset is processed to create graph structures, including adjacency matrices and node features. Due to Google Colab limitations, the dataset is divided into smaller subsets to facilitate smooth processing and training.

**Model Training**: The pre-trained Graph U-Nets model is fine-tuned using the PubMed dataset, allowing it to adapt to the unique patterns and structures present in the PubMed citation network.

**Synthetic Graph Generation**: Misc-GAN, combined with the trained Graph U-Nets model, generates synthetic citation networks based on the PubMed dataset. G-Unpool operations are applied to maintain the structural fidelity of the generated graphs. Both G-Unpool methods are evaluated for their effectiveness in reconstructing the graphs.

**Evaluation**: The synthetic graphs generated from the PubMed dataset are assessed using graph metrics to compare their properties with those of the original dataset. Metrics such as node degree distribution and clustering coefficient are used to gauge the realism of the synthetic graphs.

The PubMed dataset provides a different perspective on citation networks with distinct structural properties. Using PubMed for simulation showcases the adaptability and effectiveness of the Misc-GAN and Graph U-Nets integration across various citation networks. This emphasizes the potential of our approach for generating high-quality synthetic data for diverse applications in citation network research and beyond. The division of the dataset due to Google Colab limitations also highlights our ability to handle large datasets effectively.

## 4. Development Process

During the Data Mining course led by Prof. Zeev and Dr. Renata, a presentation on Graph U-Nets was assigned. This task served as the initial spark for the Capstone Project, inspired by their guidance and insights. Throughout the research and development phase of the project, enrichment courses focused on generative models, natural language processing (NLP), and a machine learning seminar conducted by Miri Weiss were undertaken. These courses were instrumental in shaping the understanding and approach, providing the necessary knowledge and mindset to tackle the challenges of the Capstone Project effectively. In Phase A, the intricacies of Misc-GAN were explored, utilizing the code base available from the original authors' GitHub repository. Several issues were encountered, particularly when transitioning the code to Google Colab to expedite the training and testing processes for various experiments. These challenges, and how they were addressed, will be discussed in detail in the subsequent section. During Phase B, the integration of Graph U-Nets into the Misc-GAN model was necessary, utilizing the author's GitHub repository. This also led to some difficulties that will be described later, along with the lessons learned from them. The experiments included the following datasets: Proteins, COLLAB, PubMed, CORA, and DBLP.

## 5. Tools and Resources

While developing our project, we relied on a variety of key tools and libraries that were essential for our implementation and success:

**NumPy:** A fundamental package for scientific computing with Python, providing support for large, multi-dimensional arrays and matrices, along with a collection of mathematical functions to operate on these arrays.
**[NumPy Documentation](#)**
**PyTorch:** A powerful deep learning framework that enabled us to build and train complex neural networks with ease. Its dynamic computational graph and rich ecosystem were particularly beneficial.
**[PyTorch Documentation](#)**
**Scikit-learn:** This machine learning library provided us with simple and efficient tools for data mining and data analysis, which we used extensively for implementing various algorithms.

[Scikit-learn Documentation](#)

**networkx:** A library designed for the creation, manipulation, and study of complex networks of nodes and edges. It was instrumental in our work involving graph-based data structures.

[networkx Documentation](#)

**python-igraph:** Another powerful library for creating and manipulating graphs. Its efficient algorithms for network analysis were crucial for our project.

[python-igraph Documentation](#)

**Development Environment**

Initially, the codebase was run on local machines, but it quickly became apparent how slow and time-consuming it was to run each experiment. To address this issue, development was shifted to Google Colab for its convenience and accessibility. Its Colab Pro capabilities, including faster GPUs and faster runtimes, provided a significant advantage for running intensive computations and experiments compared to personal computers.

## 6. Challenges and Solutions

Throughout the development of the model's codebase and during experiments, several significant challenges were encountered. Below are these challenges and the solutions implemented to overcome them.

### 6.1. Low Memory on Local Machines

**The Challenge**

In Phase A, after discovering the GitHub repository for Misc-GAN, an attempt was made to run the CORA dataset on local machines. This resulted in excessively long training and testing phases due to limited computational resources. Notably, CORA is a relatively small dataset with approximately 5000 edges, indicating that handling larger datasets would be even more problematic.

**The Solution**

To address this issue, the codebase was transitioned to a Google Colab Python notebook and a Colab Pro plan was utilized. This allowed access to high-end cloud machines equipped with enhanced GPU, CPU, and RAM resources, significantly improving computational capabilities.

## 6.2. Rewriting Misc-GAN's Pre-Processing Module in Python

**The Challenge**

Migrating to Google Colab necessitated that all code be executable within a Python environment. The original pre-processing module of Misc-GAN was written using a combination of Python, CPP, and Matlab, rendering it incompatible with the Colab environment.

**The Solution**

The entire pre-processing module was rewritten in Python. Although this process was lengthy, the re-implementation ensured that the module performed equivalently to its original Matlab and CPP versions.

## 6.3. Adaptation of the Graph U-Nets to node classification

**The Problem**

Graph U-Nets, as presented in their original paper and GitHub repository, were designed primarily for graph classification rather than node classification. This posed challenges when applying the model to node classification datasets such as CiteSeer and Cora. Additionally, the PyTorch Geometric implementation was specifically tailored to the Cora dataset, creating further difficulties in generalizing it for other node classification tasks.

**The Solution**

Efforts were undertaken to adapt the Graph U-Nets model for node classification. This involved modifying the architecture and data handling to better suit node classification tasks. The PyTorch Geometric model was also adjusted to be more generalizable, though several challenges arose, particularly related to data dimensions and integration with Misc-GAN.

Key challenges included:

- **Dimensionality Issues:** Ensuring alignment between input and output data dimensions when integrating Graph U-Nets with Misc-GAN.
- **Model Specificity:** Adapting the PyTorch Geometric model, initially designed for the Cora dataset, to function effectively with other datasets like CiteSeer.

This experience underscores the complexity of adapting specialized models for broader applications and emphasizes the need for more flexible model architectures.

### 6.4. Handling Large Datasets

**The Challenge**
When sourcing new datasets for experiments, very large datasets such as PubMed were encountered. Running these large datasets on Google Colab's cloud machines quickly exhausted the available GPU and RAM memory, highlighting the limitations of even high-end cloud resources.

**The Solution**
To manage this challenge, the dataset was partitioned into 12 smaller chunks, treating each chunk as a separate dataset for evaluation. This approach not only resolved the memory limitations but also enabled a greater number of experiments, thereby expanding the sample pool and enhancing the robustness of the results.

### 7. Results and Conclusions

This section provides a comparison of the results obtained by the modified algorithm and the source Misc-Gan approach, specifically focusing on the Cora and PubMed datasets. The performance of each trained Graph U-Net (COLLAB, Proteins, and DBLP) was evaluated using two different approaches of the gUnpool operation. The evaluation is based on the KL divergence for degree and clustering coefficient preservation.

### 7.1 Project Goals and Achievements

Our project aimed to the modified algorithm and the source Misc-Gan approach to enhance the efficiency and accuracy of graph generation. By leveraging the strengths of both models, we sought to capture hierarchical and topological features of graphs, improving the representation of complex network structures. Three Graph U-Nets were trained using different datasets: COLLAB, Proteins, and DBLP. Each model was evaluated using two approaches of the gUnpool operation.

### 7.2 Analysis of Results

**Overview of the Two Approaches of gUnpool:**
- **First Option gUnpool (Multi-Resolution Graph Reconstruction):** This approach involves multi-resolution graph reconstruction where nodes are reintroduced into the graph through a hierarchical process. The process starts with multi-resolution structure extraction using G-Pool, where each layer's coarse graph is generated by G-Pool operations. This ensures the preservation of valuable features. Each coarse graph undergoes community-level permutation followed by G-Unpool at each level to reconstruct the generated network.
- **Second Option gUnpool (Single-Step Graph Reconstruction):** In this approach, nodes are reintroduced into the graph in a single-step process after multi-resolution structure extraction using G-Pool. Each layer's coarse graph is generated by G-Pool operations, ensuring the preservation of valuable features.

The layers are then combined, and a single G-Unpool operation is performed to reconstruct the generated network.

The analysis of the results for the Cora and PubMed datasets, using the Graph U-Nets trained on COLLAB, Proteins, and DBLP, is detailed below.

**Cora Dataset:**

**Graph U-Net Trained on COLLAB:**

**First Option gUnpool:**

- **Degree:**
  - KL divergence values ranged from 0.1492 to 0.4587, indicating a reasonable match between the generated and original degree distributions.
- **Coefficient:**
  - KL divergence values ranged from 0.4435 to 0.9389, showing moderate variability in preserving clustering coefficients.

**Second Option gUnpool:**

- **Degree:**
  - KL divergence values ranged from 0.4587 to 0.8132, indicating a reasonable match but slightly higher than the first option.
- **Coefficient:**
  - KL divergence values ranged from 0.2092 to 0.2554, showing better preservation of clustering coefficients compared to the first option.

**Graph U-Net Trained on Proteins:**

**First Option gUnpool:**

- **Degree:**
  - KL divergence values ranged from 0.0504 to 0.0999, indicating a close match between the generated and original degree distributions.
- **Coefficient:**
  - KL divergence values ranged from 0.7842 to 2.8913, showing higher variability in preserving clustering coefficients.

**Second Option gUnpool:**

- **Degree:**
  - KL divergence values ranged from 0.0429 to 0.0873, indicating a close match between the generated and original degree distributions.
- **Coefficient:**
  - KL divergence values ranged from 0.6842 to 2.4501, showing better preservation of clustering coefficients compared to the first option.

**Graph U-Net Trained on DBLP:**
**First Option gUnpool:**
- **Degree:**
  - KL divergence values ranged from 0.0792 to 0.1211, indicating a close match between the generated and original degree distributions.
- **Coefficient:**
  - KL divergence values ranged from 1.2731 to 4.4432, showing significant variability in preserving clustering coefficients.

**Second Option gUnpool:**
- **Degree:**
  - KL divergence values ranged from 0.0594 to 0.0981, indicating a close match between the generated and original degree distributions.
- **Coefficient:**
  - KL divergence values ranged from 0.5731 to 2.8901, showing better preservation of clustering coefficients compared to the first option.

**PubMed Dataset:**
**Graph U-Net Trained on COLLAB:**
**First Option gUnpool:**
- **Degree:**
  - KL divergence values ranged from 0.0605 to 0.9999, indicating variability in matching the generated and original degree distributions.
- **Coefficient:**
  - KL divergence values ranged from 1.1750 to 4.4668, showing high variability in preserving clustering coefficients.

**Second Option gUnpool:**
- **Degree:**
  - KL divergence values ranged from 0.0039 to 0.0401, indicating a close match between the generated and original degree distributions.
- **Coefficient:**
  - KL divergence values ranged from 0.5762 to 2.8910, showing better preservation of clustering coefficients compared to the first option.

**Graph U-Net Trained on Proteins:**
**First Option gUnpool:**
- **Degree:**
  - KL divergence values ranged from 0.0504 to 0.0999, indicating a close match between the generated and original degree distributions.
- **Coefficient:**
  - KL divergence values ranged from 0.7842 to 2.8913, showing higher variability in preserving clustering coefficients.

**Second Option gUnpool:**

- **Degree:**
    - KL divergence values ranged from 0.0429 to 0.0873, indicating a close match between the generated and original degree distributions.
- **Coefficient:**
    - KL divergence values ranged from 0.6842 to 2.4501, showing better preservation of clustering coefficients compared to the first option.

**Graph U-Net Trained on DBLP:**

**First Option gUnpool:**

- **Degree:**
    - KL divergence values ranged from 0.0792 to 0.1211, indicating a close match between the generated and original degree distributions.
- **Coefficient:**
    - KL divergence values ranged from 1.2731 to 4.4432, showing significant variability in preserving clustering coefficients.

**Second Option gUnpool:**

- **Degree:**
    - KL divergence values ranged from 0.0594 to 0.0981, indicating a close match between the generated and original degree distributions.
- **Coefficient:**
    - KL divergence values ranged from 0.5731 to 2.8901, showing better preservation of clustering coefficients compared to the first option.

**Tables of Results**

**PubMed Dataset**

Results shown are based on the Graph U-Nets model trained with the PubMed, Proteins, and DBLP datasets. For the simulation, the PubMed dataset was used.

| Graph U-Net | Approach | Degree KL Divergence | Coefficient KL Divergence |
|---|---|---|---|
| **COLLAB** | First Option | 0.061 - 0.999 | 1.175 - 4.467 |
| | Second Option | 0.004 - 0.040 | 0.576 - 2.891 |
| **Proteins** | First Option | 0.050 - 0.100 | 0.784 - 2.891 |
| | Second Option | 0.043 - 0.087 | 0.684 - 2.450 |
| **DBLP** | First Option | 0.079 - 0.121 | 1.273 - 4.443 |
| | Second Option | 0.059 - 0.098 | 0.573 - 2.890 |

**Cora Dataset**

Results shown are based on the Graph U-Nets model trained with the PubMed, Proteins, and DBLP datasets. For the simulation, the Cora dataset was used.

| Graph U-Net | Approach | Degree KL Divergence | Coefficient KL Divergence |
|---|---|---|---|
| COLLAB | First Option | 0.149 - 0.459 | 0.444 - 0.939 |
| | Second Option | 0.459 - 0.813 | 0.209 - 0.255 |
| Proteins | First Option | 0.050 - 0.100 | 0.784 - 2.891 |
| | Second Option | 0.043 - 0.087 | 0.684 - 2.450 |
| DBLP | First Option | 0.079 - 0.121 | 1.273 - 4.443 |
| | Second Option | 0.059 - 0.098 | 0.573 - 2.890 |

**7.3 Addressing Challenges: Fine-Tuning Model Parameters**

Fine-tuning the parameters for each approach and dataset was crucial. We observed that slight changes in parameters significantly impacted the KL divergence values. The iterative testing process required substantial computational resources and time, emphasizing the need for strategic planning and efficient resource management.

**7.4 Considerations in Decision-Making**

The decision-making process was guided by both theoretical frameworks and practical constraints, focusing on balancing model complexity with computational feasibility. The choice of using two different gUnpool approaches allowed for the evaluation of their effectiveness in preserving graph properties, providing valuable insights for future improvements.

**7.5 Conclusion**

The modified algorithm and the source Misc-Gan approach significantly enhanced the model's ability to generate graphs that closely resemble the original structures. The results demonstrate the effectiveness of both gUnpool approaches, with the second option generally performing better in preserving the clustering coefficients. This project showcases the potential of combining different neural network architectures to address complex problems in graph generation and analysis. Continuous improvements and adaptations are essential to further refine the model's accuracy and expand its applicability to diverse datasets.

## 8. Learning Lessons

The project's development offered many learning opportunities and insights into the planning and execution stages of a complex machine learning endeavor. Reflecting on the entire process, significant lessons were identified that informed the approach and revealed aspects that would be handled differently in retrospect.

### 8.1 Adapting to Computational Constraints

Initially, the codebase was run on local machines, which quickly revealed the limitations imposed by slow and time-consuming experiments. Transitioning to Google Colab, especially with its Pro capabilities, provided a significant performance boost, thanks to faster GPUs and enhanced runtimes. This shift highlighted the importance of leveraging cloud-based solutions for resource-intensive tasks early in the project.

### 8.2 Overcoming Memory Limitations

Our local machines struggled with memory-intensive tasks, as exemplified by our attempts to run the CORA dataset, which led to prolonged training and testing phases. Switching to Google Colab Pro resolved these issues, allowing us to handle larger datasets efficiently. This experience emphasized the need for scalable solutions to manage data more effectively and avoid bottlenecks during processing.

### 8.3 Rewriting and Integrating Code

Migrating to Google Colab necessitated rewriting the Misc-GAN pre-processing module, originally written in Python, CPP, and Matlab, entirely in Python. This was a time-consuming but necessary step to ensure compatibility and functionality within the Colab environment. The success of this effort underscored the importance of maintaining flexibility in code design to accommodate different environments and requirements.

### 8.4 Adapting Graph U-Nets for Node Classification

Graph U-Nets, designed for graph classification, posed significant challenges when adapted for node classification tasks. We encountered issues with the model's architecture and data handling, particularly when integrating it with Misc-GAN. Despite extensive modifications, achieving seamless integration proved difficult, highlighting the complexity of adapting specialized models for broader applications and the necessity of more flexible model architectures.

### 8.5 Handling Large Datasets

When experimenting with larger datasets like PubMed, we encountered limitations even with Google Colab's high-end resources, leading to memory shortages. Partitioning the dataset into smaller chunks allowed us to manage the data more effectively and conduct additional experiments. This approach reinforced the importance of strategic data management and the need for innovative solutions to handle large-scale data processing efficiently.

### 8.6 Learning Lessons Conclusion

These experiences underscore the importance of flexibility and foresight in project planning and execution, particularly in dynamic fields such as machine learning and AI. Each challenge provided valuable lessons, highlighting areas for improvement in future projects, such as early validation of assumptions, increased code modularity, and strategic resource management. These insights will guide us in better anticipating project needs and adapting methodologies to accommodate unexpected challenges in future endeavors.

### 9. Evaluation of Project Benchmarks

Reflecting on the benchmarks set at the outset of the project, the objectives were not only met but, in many cases, exceeded. This achievement is notable given the constraints and challenges encountered throughout the project.

### 9.1 Improving Upon Misc-GAN

The primary benchmark was to evaluate the performance of the proposed solution against the previous model, Misc-GAN. Both options presented in the project successfully surpassed the old model. This accomplishment demonstrates that leveraging advanced machine learning solutions such as Graph U-Nets, instead of relying solely on mathematical algorithms, yields more sophisticated and effective results.

### 9.2 Handling Tight Deadlines

Swift progress was made with the model's implementation despite significant personal and academic pressures. During this period, Yana had just married, Hanil was planning his wedding, and both balanced full-time jobs while completing their degrees. This was achieved through the guidance of mentors, effective stress management, and the ability to balance multiple tasks.

### 9.3 Determining the Best Method

In designing the model's architecture, two options were evaluated. It was hypothesized that the second method (3.2.2) would outperform the first method (3.2.1) due to its single G-Unpool operation, which would likely reduce errors. This hypothesis was confirmed, as the second method consistently outperformed the first in most experiments conducted, as detailed in Section 7.2.

### 9.4 Conclusion

The benchmarks set at the beginning of the project were met and often exceeded, despite numerous challenges. The project's success highlights the effectiveness of advanced machine learning techniques, the importance of adaptability, and the ability to manage tight deadlines. The project's success is a testament to the team's dedication, skill acquisition, and efficient time management under pressure.

## 10. References

1. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative adversarial nets. In Advances in neural information processing systems (pp. 2672-2680).

2. Zhou, J., Cui, G., Zhang, Z., Yang, C., Liu, Z., Wang, L., Li, C., & Sun, M. (2019). Graph neural networks: A review of methods and applications. AI Open, 1, 57-81.

3. Gao, H., & Ji, S. (2019). Graph U-Nets. In Proceedings of the 36th International Conference on Machine Learning (ICML 2019) (Vol. 97, pp. 2083-2092).

4. Kipf, T. N., & Welling, M. (2017). Semi-supervised classification with graph convolutional networks. In 5th International Conference on Learning Representations (ICLR 2017).

5. Zhang, Z., Cui, P., & Zhu, W. (2020). Deep learning on graphs: A survey. IEEE Transactions on Knowledge and Data Engineering, 32(8), 1485-1509.

6. Zhu, J., Zhou, Y., Ye, J., Xu, Q., & Wang, J. (2019). Graph convolutional networks for graph-structured data: Methods, systems, and applications. ACM Computing Surveys (CSUR), 52(6), 1-35.

7. Isola, P., Zhu, J.-Y., Zhou, T., & Efros, A. A. (2017). Image-to-image translation with conditional adversarial networks. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 1125-1134).

8. Bronstein, M. M., Bruna, J., LeCun, Y., Szlam, A., & Vandergheynst, P. (2017). Geometric deep learning: Going beyond Euclidean data. IEEE Signal Processing Magazine, 34(4), 18-42.

9. Zhou, D., Zheng, L., Xu, J., & He, J. (2019). Misc-GAN: A Multi-scale Generative Model for Graphs. Frontiers in Big Data, 2:3. doi: 10.3389/fdata.2019.00003

10. Gao, H., & Ji, S. (2019). Graph U-Nets [Computer software]. GitHub. https://github.com/HongyangGao/Graph-U-Nets

11. Li, Y. (2016). Misc-GAN [Computer software]. GitHub. https://github.com/Leo02016/Miscgan