# Anomalies Detection Using Graph Diffusion Model

**Students:**

Roei Nizri – RoeiNizri1@gmail.com, 318260742

Ran Polac – RanPolac@gmail.com, 318265808


**Supervisors:** Dr. Renata Avros, Prof. Zeev Volkovich

# Table of Contents

# Abstract

Anomaly detection on graphs presents significant challenges, primarily due to label imbalance and data insufficiency. These factors make learning and finding real anomalies in complex graph data hard. To address this, the study introduces a diffusion model-based graph generator capable of simultaneously producing graph structures and node features in a multitask fashion, focusing on generating synthetic anomalous examples. By integrating these synthetic anomalies into the training dataset, we achieve a more balanced dataset, significantly enhancing the training and performance of downstream anomaly detection models. The approach addresses the challenges of label imbalance and data insufficiency in graph anomaly detection, providing a robust framework for generating synthetic anomalies and improving the effectiveness of anomaly detection models.

## 1. Introduction and Background

Anomaly detection has been a critical task in the realm of machine learning and data mining applications, due to its significance in identifying irregular patterns that deviate from the norm in large datasets. The main issue is the difficulty in identifying anomalous nodes within graph-structured data due to the scarcity of anomaly labels compared to normal data. This imbalance makes it hard for detection models to learn effectively, leading to suboptimal performance in detecting true anomalies. Several approaches have been proposed for anomaly detection on graphs. Methods like Graph Convolutional Networks [2] (GCNs) and Graph Attention Networks (GATs) [4] learn node embeddings and identify nodes that deviate from learned patterns. However, these methods often struggle with label imbalance, making it difficult to identify rare anomalies accurately. Random Walk Algorithms use random walks to capture proximity and structural patterns in the graph, identifying anomalies based on deviations in the random walk distributions. However, they may be computationally expensive and less effective in large or highly connected graphs. Factorizing the adjacency matrix of the graph to uncover latent features, with anomalies detected based on outliers in the latent space. This method may not fully capture the intricate relationships and dependencies within the graph, leading to suboptimal anomaly detection. To address the limitations of existing methods, we propose a diffusion model-based graph generator. This generator is capable of simultaneously producing graph structures and node features in a multitask fashion, with a particular focus on generating synthetic anomalous examples. By integrating these synthetic anomalies into the training dataset, we achieve a more balanced dataset, significantly improving the training and performance of downstream anomaly detection models. The solution is designed to benefit the following stakeholders.

1. **Data Scientists and Machine Learning Engineers：** These professionals are constantly seeking advanced techniques to improve the accuracy and efficiency of anomaly detection systems.

2. **Cybersecurity Experts:** Cybersecurity professionals use anomaly detection to identify suspicious activities and potential threats within network traffic or user behavior.

3. **Financial Analysts and Fraud Detection Specialists:** These individuals are interested in detecting fraudulent activities, such as abnormal transaction patterns or insider trading, which can be represented as graphs.

The proposed model enhances the detection of fraudulent transactions and anomalies in graph-structured data by generating synthetic anomalies and addressing label imbalance. This approach leads to more accurate and reliable data, improving overall security measures. In the following sections, we will delve deeper into the methodology of the diffusion model-based graph generator, discuss its advantages over existing methods, and present experimental results demonstrating its efficacy in anomaly detection tasks.

## 2. Related Work (Literature Review)

Anomaly detection on graphs is growing quickly because of the unique challenges with graph data. Unlike regular tabular data, graphs show complex relationships [13] and dependencies between entities, making anomaly detection harder. Early methods, such as clustering and proximity-based techniques, were designed to handle simpler graph structures but struggled with more complex, dynamic graphs. For example, proximity-based methods identify anomalies by measuring the similarity between nodes, which work well for small, static graphs but fail in larger, evolving graphs with diverse relationships. To address these challenges, researchers have developed various algorithms that leverage the structural properties of graphs. **Graph-based Clustering Algorithm** [1] algorithms segment the graph into clusters based on node similarity. Nodes that do not fit well into any cluster are considered anomalies. However, these methods struggle with high-dimensional and dynamic graphs. **Graph autoencoders** [2] learn low-dimensional representations of nodes and reconstruct the graph structure. Anomalies are identified as nodes with high reconstruction errors. While effective, autoencoders require careful tuning and may not generalize well to different types of graphs.
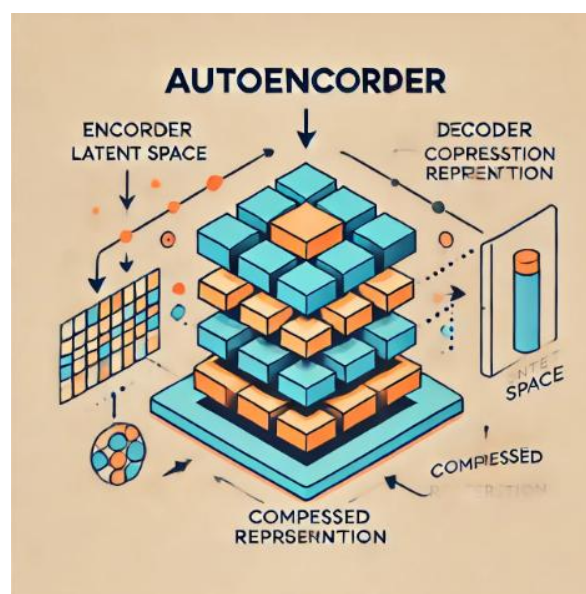


Figure 1: The architecture of a Autoencoder

**Generative Adversarial Networks** [3] create synthetic graph data to improve the training set and fix label imbalance. Although promising, GANs [11] require a lot of computing power and can be hard to train.
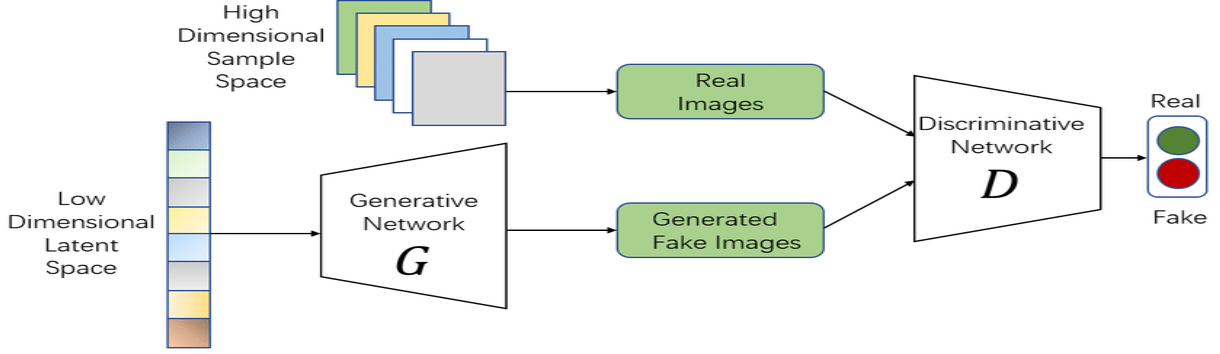


Figure 2: The architecture of a basic GAN

## 2.1 Graph Neural Networks (GNNs)

Graph Neural Networks (GNNs) [1] have gained prominence due to their ability to learn complex patterns and dependencies within graph data. GNNs can capture both local and global structures, making them effective for detecting anomalies in various applications such as fraud detection, network security, and social network analysis. One notable technique is the Graph Convolutional Network (GCN) [17], which extends the concept of convolutional neural networks to graph-structured data. GCNs aggregate information from neighboring nodes, allowing them to learn node representations that incorporate the graph's structure. This method has been successful in identifying anomalies in social networks and citation networks where the relationships between entities are crucial for detecting irregular patterns. However, GCNs often require large, labeled datasets for training, which can be a limitation in scenarios with scarce anomaly labels. Another advanced GNN technique is the Graph Attention Network (GAT), which introduces an attention mechanism to weigh the importance of different neighbors during the aggregation process. This allows GATs to focus on the most relevant parts of the graph, improving the detection of subtle anomalies. Despite their effectiveness, GATs can be computationally intensive and may still struggle with highly imbalanced datasets.
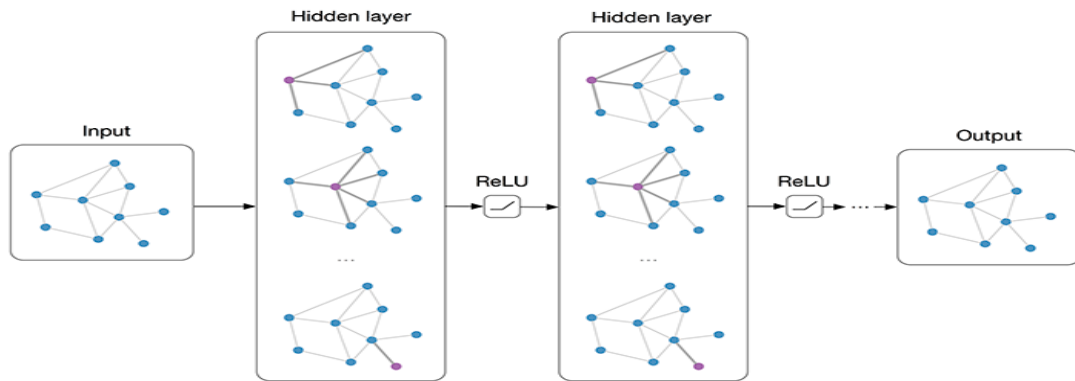


Figure 3: Multi-layer Graph Convolutional Network (GCN) with first-order filters

[5]

## 2.2 Encoding with Graph Autoencoder

Graph autoencoders [2] are used to map anomaly detection graphs into a latent space, simplifying the graph data while retaining its most salient features. The autoencoder consists of two main parts:

**Encoder:** Transforms graph data into a lower-dimensional latent space by compressing the data.

**Decoder:** Reconstructs the input data from the latent representation, ensuring that the learned representations are meaningful.

Graph autoencoders [2] have been effective in identifying anomalies as nodes with high reconstruction errors. However, they require careful tuning and may not generalize well to different types of graphs. Variational Graph Autoencoders (VGAEs) [4] [15] extend this approach by incorporating a probabilistic framework, which improves generalization by predicting a distribution over the latent space. This allows VGAEs to better handle variability in graph data and generate more realistic reconstructions.

## 2.3 Generating Graph Structure

A traditional autoencoder is first applied to generate node attributes. This autoencoder learns efficient representations of the data, capturing the essential characteristics while reducing dimensionality. The process involves compressing the input data through the encoder and reconstructing it via the decoder.

**Training:** The model adjusts weights to minimize the difference between the original input and its reconstruction using a loss function called Mean Squared Error (MSE). Through this process, the autoencoder learns to prioritize the most significant features in the input data, effectively learning a compressed but informative representation. Variational Autoencoder (VAE): Enhanced Autoencoder: Introduces a probabilistic approach to improve generalization.

**Process:** The encoder predicts two parameters ($\mu$, $\mu$) and standard deviation ($\sigma$) for a Gaussian distribution. During the forward pass, a sample is drawn from this distribution and decoded to reconstruct the input data.

**Loss Function:** Combines the reconstruction term (MSE) between input and reconstructed output) and a regularization term (Kullback-Leibler (KL)) [9] divergence between the learned distribution and a standard Gaussian). The VAE improves by predicting Gaussian distribution parameters (mean and standard deviation) and using a combined reconstruction and regularization loss function.

Conditional Generation for Nodes with Positive Labels controls over generating different node types (anomalous and normal). The methodology for generating graph structures employs the Variational Graph Autoencoder (VGAE), which captures both topological and feature information into a latent space. The VGAE model uses a shared Graph Neural Network (GNN) encoder for multi-task learning, adapting the VAE architecture to graph data. To handle heterogeneous graphs with multiple node and edge types, the standard VGAE encoder is replaced with a Heterogeneous Graph Transformer (HGT) [7]. The HGT can process diverse node and edge types, allowing for the accurate generation of adjacent matrices specific to each edge type. This approach enhances the model's ability to capture complex relationships in heterogeneous graphs, making it more effective for anomaly detection in such data. Different decoders are used for different edge types, allowing accurate generation of adjacency matrices specific to each edge. The transition to matrices allows a more accurate comparison between the weights of the edges and nodes.

## 2.4 Timestamp Generation

A temporal generator designed to generate timestamps parallels the structure of the feature generator but targets singular dimensional output. The generated timestamps are integrated into the diffusion model's latent space, ensuring that the synthetic data generated (including anomalies) maintains realistic temporal characteristics. Timestamps help capture the evolution of the graph over time, enabling the detection of anomalies that occur due to changes in the graph structure or node attributes. This temporal aspect is crucial for applications like network intrusion detection and fraud detection, where the timing of events can provide important clues about anomalous behavior.

**Purpose:** Generate timestamps as part of the graph data.

**Architecture:** Like the feature generator but focused on generating a single continuous variable representing time.

**Loss Function:** Uses MSE to minimize the average squared difference between estimated and actual timestamps, ensuring precision and reliability in the generated temporal data.

## 2.5 Diffusion models

Diffusion models, also called diffusion probabilistic models or score-based generative models, gradually transform simple distributions, like Gaussian noise, into the desired data distribution. This transformation is achieved through a series of small, reversible steps, allowing the model to learn the underlying data structure effectively. Recent advances in generative models, especially diffusion models, offer a promising approach to addressing the challenges of anomaly detection on graphs. Initially developed for image synthesis, these models have shown exceptional capability in capturing complex data distributions through a process of gradual noising and denoising. Although their application to graph data is relatively new, it holds significant potential. One key advantage is their ability to generate realistic graph structures and node features, which can be used to create synthetic anomalies. Consequently, this approach helps mitigate the label imbalance problem by enriching the training datasets with synthetic anomaly nodes. By generating synthetic graphs that mirror the intricacies of real-world graphs, these models provide a robust foundation for training anomaly detection algorithms. Moreover, they excel at capturing the complex, high-dimensional relationships inherent in graph data, making them well-suited for generating realistic synthetic anomalies. By learning the normal patterns within the graph data, the model can effectively identify anomalies as deviations from these learned patterns. This process involves adapting the diffusion process to work with the complex relationships and structures inherent in graph data. In the context of anomaly detection, the Denoising Diffusion Probabilistic Model (DDPM) [12] enhances the quality of synthetic anomaly generation. DDPM operates by adding noise to data points and then learning to reverse this process, effectively generating data points from noise. This process involves adapting the diffusion process to work with the complex relationships and structures inherent in graph data. In summary, by leveraging their unique capabilities, researchers can create enriched datasets that improve the robustness and accuracy of anomaly detection algorithms, making these models a powerful tool in the ongoing effort to enhance anomaly detection on graphs.
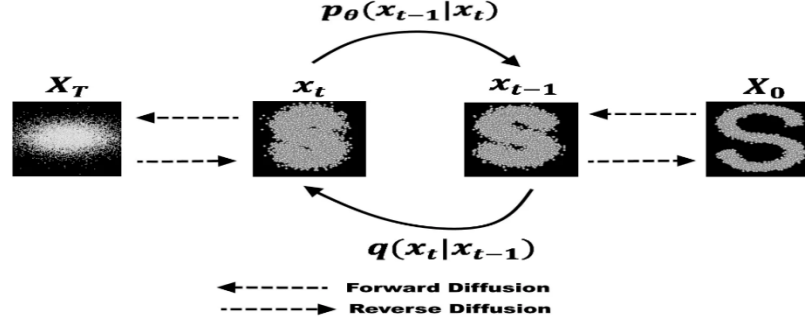
Figure 4: The main idea of diffusion model. By adding noise in the diffusion process the image became more cleaner and by imply reverse diffusion process the image return to the origin.

## 2.5.1 Diffusion Models in Latent space

Diffusion models in latent space offer a sophisticated approach to generating synthetic data, including anomalies, for anomaly detection tasks. These models operate by transforming input data into a latent representation, adding controlled noise through a diffusion process, and then applying a denoising process to reconstruct the data. This methodology allows the generation of realistic graph structures and node features that closely resemble real-world data while incorporating synthetic anomalies essential for robust training of anomaly detection models. The process begins with encoding the input graph data, which includes graph nodes and their features, into a latent space. The encoding phase utilizes techniques such as autoencoders or variational autoencoders (VAEs) to compress the high-dimensional graph data into a lower-dimensional latent representation. This latent representation retains the most significant features and structures of the original data, making it more manageable and suitable for further processing. Once the data is encoded into the latent space, the diffusion process introduces noise incrementally. The diffusion process can be mathematically described as a series of steps where Gaussian noise is added to the latent representation based on a predefined variance schedule. This gradual introduction of noise transforms the data into a more chaotic state, which the model learns to reverse. The key idea is to train the model to understand the underlying data distribution by learning how to effectively add and remove noise. The denoising process follows, where the model learns to reverse the added noise, reconstructing the original data from the noisy latent representations. This process involves training a neural network to predict the clean data from its noisy version at each step. By iteratively applying this denoising process, the model generates data that closely mirrors the original input, including the generation of synthetic anomalies.
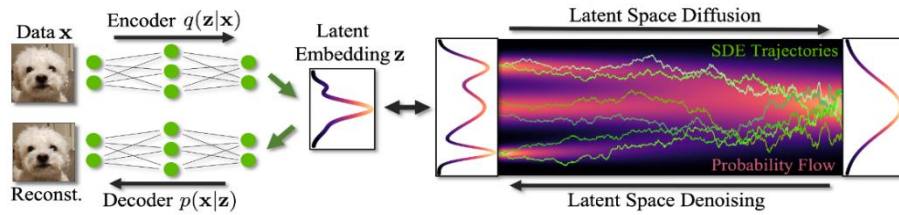


Figure 5: The process of the model, the data (citations) are encoded into a latent space, where a diffusion process introduces controlled noise, generating synthetic anomalies. The denoised data is then decoded back, allowing for robust anomaly detection.

## 2.6 DeepWalk and Node2Vec

DeepWalk [5] and Node2Vec [6] are algorithms that generate node embeddings by simulating random walks on the graph. These embeddings capture the structural relationships between nodes, allowing for the identification of anomalies as nodes with embeddings that deviate significantly from the norm. DeepWalk performs truncated random walks on the graph and treats these walks as sentences to train a Skip-Gram model, like word2vec in natural language processing. Node2Vec extends this by using a biased random walk that balances between breadth-first and depth-first strategies, allowing for more flexible graph exploration. While effective for homogeneous graphs, these methods struggle with heterogeneous graphs and may miss complex relationships.

## 2.7 Generative Adversarial Networks (GANs)

Generative Adversarial Networks (GANs) [3] generate synthetic graph data to augment the training set and address label imbalance. GANs consist of two neural networks: a generator that creates synthetic data and a discriminator that evaluates the authenticity of the data. [11] This adversarial process drives the generator to produce increasingly realistic data. While GANs have shown promise in generating high-quality synthetic data, they can be computationally intensive and challenging to train. The training process requires careful balancing between the generator and discriminator to prevent issues like mode collapse. Despite these challenges, GANs offer a powerful tool for generating synthetic anomalies and enriching the training datasets for anomaly detection models. In conclusion, the field of anomaly detection on graphs has seen significant advancements with the development of machine learning and generative models. Graph Neural Networks [1] and generative models, particularly diffusion models, offer powerful tools to address the unique challenges of graph-structured data. The proposed diffusion model-based graph generator represents a approach to mitigate label imbalance and enhances the training of anomaly detection algorithms. As the field continues to evolve, these advancements hold great promise for improving the detection of anomalies in complex, dynamic graph data.
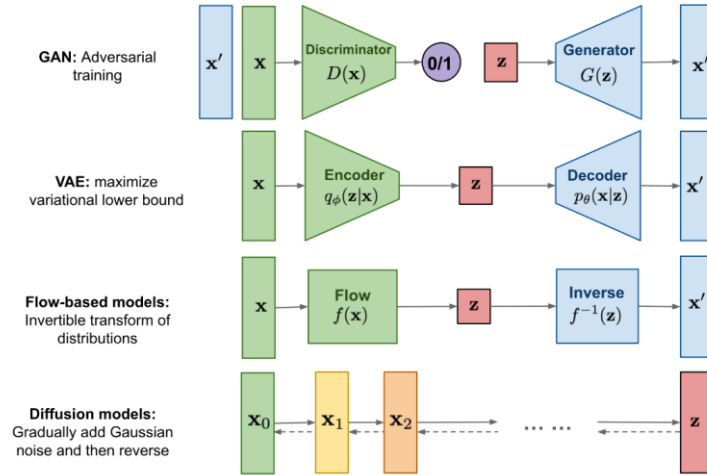


Figure 6: Overview of different types of generative models.

# 3. Engineering Process

1. **Prepare Dataset**
Gather a large dataset of citations, including metadata such as authors, publication dates, and sources. Represent the dataset as a graph where: nodes represent individual citations. Edges represent relationships between citations, such as being from the same author, same publication, or similar content.
*Input:* A collection of citations with associated metadata.
*Output:* A graph representation of the dataset with nodes as citations and edges as relationships between citations.

2. **Graph Construction**
Each node (quote) should have features such as text embeddings, author information, and publication details. Create edges based on relationships like shared authorship, publication proximity, or semantic similarity.
*Input:* Dataset with node features and edge criteria.
*Output:* A fully constructed graph with nodes and edges defined.

3. **Graph Autoencoder implementation**
We use a graph neural network (GNN) to encode the graph into a latent space, capturing the structural and feature information. Train the autoencoder to reconstruct the original graph from the latent representation.
*Input:* Constructed graph with node features and edges.
*Output:* Latent space representation of the graph and reconstructed graph from the autoencoder.

4. **Diffusion Model application**
Implement a diffusion model in the latent space to generate synthetic anomalies. This involves iterative noise addition and removal to create synthetic nodes (citations) that resemble anomalies. We ensure the model can generate nodes with specific labels (anomalous or normal) based on the citation's characteristics.
*Input:* Latent space representation of the graph.
*Output:* Synthetic anomalous nodes in the latent space and their corresponding graph representations.

5. **Anomaly Detection**
We use synthetic anomalies along with real data to train an anomaly detection model. Incorporate both real and synthetic data to address label imbalance.
Assess the model's performance using metrics like AUROC and AUPRC to ensure it effectively identifies anomalies.
*Input:* Real and synthetic graph data with labels.
*Output:* Anomaly detection model capable of identifying anomalous citations with performance metrics.

## 3.1 Pre-processing

The purpose goal is to enrich the dataset with synthetic anomaly nodes for robust training of anomaly detection systems. These synthetic nodes include four key components:
**Node Feature:** Represents various attributes associated with the node, such as metrics or characteristics that describe its properties and behavior within the graph.
**Node Timestamp:** A temporal marker indicating when a specific event or interaction involving the node occurred, crucial for analyzing the temporal dynamics of the network.
**Positive Label:** A binary indicator denoting the anomaly nature of the node, used to differentiate between normal and anomalous behavior for training and validating detection models.
**Heterogeneous Graph Structure:** Different types of nodes and edge connections.

### 3.1.1 Graph Autoencoder

The encoding process involves using a graph autoencoder to map anomaly detection graphs into a latent space. The autoencoder consists of two main parts:

**Encoder:** In the implementation for identifying suspected citations in an article, the encoder processes each citation (node) and its relationships (edges) by:

*Embedding Text Features,* converting the textual content of citations into dense vector representations using techniques like word embeddings or transformer models.

*Incorporating Metadata,* adding author information, publication details, and other relevant metadata to the node features. Using graph convolutional layers to capture the relationships between citations, like shared authorship or semantic similarity, and transforming these relationships into a compact form that keeps the important patterns and connections.

**Decoder:** Reconstructs text features and attempts to regenerate the textual content of citations from their latent representations.

*Restores Metadata and Relationships,* reconstructs the author information, publication details, and relationships between citations based on the latent space representations.

*Validates Reconstruction,* ensures that the reconstructed graph closely resembles the original input graph, maintaining the key features and structures identified by the encoder.


### 3.1.2 Generating Node Attributes

To control the creation of node types, especially anomalous ones, conditional variables are added to the VAE. During training, labels indicating node types are concatenated with feature vectors and fed into the encoder. During generation, the decoder uses these labels to create feature vectors that match the specified labels. This ensures that when the label indicates 'fraudulent' (anomalous), the generated features are characteristic of fraudulent nodes.

The benefit is to allow a more controlled generation of nodes, improving the model's ability to differentiate and generate distinct node types, which is crucial for applications like anomaly detection. In the implementation for identifying suspected citations, the process begins with preparing a graph where nodes represent citations and edges represent relationships such as shared authorship or semantic similarity. The citations have attributes including text embeddings, author information, and publication details. The autoencoder's encoder processes each citation and its relationships, compressing this information into a lower-dimensional latent space while retaining essential features. The decoder then reconstructs the original quote features from the latent representation, minimizing the reconstruction error (MSE) and ensuring the latent space captures the critical attributes of the citations. The VAE encoder predicts the mean and standard deviation for the Gaussian distribution representing each citation in the latent space. A sample is drawn from this distribution and decoded to reconstruct the citation, optimizing both the reconstruction accuracy and the regularization term (KL divergence). During training, node-type labels (normal or anomalous) are combined with feature vectors and processed by the encoder. This ensures the latent representation includes information about the node type. In the generation phase, the decoder uses these labels to produce feature vectors corresponding to the specified node type, enabling controlled generation of normal and anomalous citations. The VAE generates synthetic citations labeled as anomalous by using conditional variables. These synthetic anomalies imitate suspected citations that differ from usual patterns. The anomaly detection model is then trained on both real and synthetic data, using the synthetic anomalies to address label imbalance, enhancing the model's ability to identify suspected citations. This approach improves feature extraction, allows for precise generation of anomalous citations, and enhances the anomaly detection model's ability to detect suspected citations, ensuring robust performance.

### 3.1.3 Generating Heterogeneous Graph Structure

Following the previous step, to handle heterogeneous graphs with multiple node and edge types, the standard VGAE encoder is replaced with a Heterogeneous Graph Transformer (HGT) [16]. The HGT can process diverse node and edge types, allowing for the accurate generation of adjacency matrices specific to each edge type. This approach enhances the model's ability to capture complex relationships in heterogeneous graphs, making it more effective for anomaly detection in such data. Separate decoders are used for different edge types, allowing accurate generation of adjacency matrices specific to each edge. The transition to matrices allows a more accurate comparison between the weights of the edges and nodes. In the implementation for identifying suspected citations, the HGT begins by preparing a graph where nodes represent citations and edges represent relationships such as shared authorship or semantic similarity. The HGT processes each citation and its relationships, capturing the diverse and complex interactions within the graph. This step ensures the model can accurately generate adjacency matrices for different edge types, enhancing its ability to detect unusual patterns.

### 3.1.4 Timestamp Generation

The anomaly detection system adds synthetic anomaly nodes to datasets, using a temporal generator that works like the feature generator but focuses on a single output dimension. This temporal data significantly expands the space of time within the dataset, allowing for a more comprehensive analysis of temporal patterns and trends. By generating accurate and reliable timestamps, the system can monitor the progression and evolution of data points over time, enhancing the ability to detect anomalies that may arise due to unusual temporal behaviors. Furthermore, this step allows us to monitor the learning process of anomaly detection models. By adding accurate timestamps, we can check how well the models identify temporal anomalies. This ensures they detect anomalies based on both static features and event timing. This dual focus improves the system's ability to identify suspected citations, leading to more reliable detection outcomes.

## 3.2 Diffusion Model in Latent Space

Diffusion models are used to generate new data by adding noise to data points and then learning how to reverse this process, effectively generating data points from noise. This process can be particularly useful in generating synthetic anomalies for anomaly detection. The diffusion model in latent space is a crucial technique used to improve the generation quality of synthetic anomalies in graph data, which is essential for effective anomaly detection. In the implementation for identifying suspected citations, we adopt the Denoising Diffusion Probabilistic Model (DDPM), which stands for "Diffusion models with Denoising Priors." This type of generative model simulates the gradual diffusion (or spreading out) of noise over time. To explain this further, imagine that each citation in the dataset is a point in a complex graph. By introducing noise, we can perturb these points slightly and then train the model to learn how to revert these changes, effectively teaching it the underlying structure of normal and anomalous citations. In practice, this means we first add noise to the data points in the citation graph and then use the DDPM to denoise or reverse this process. This helps the model learn to differentiate between normal and anomalous patterns. For instance, if a citation suddenly deviates significantly from the established pattern of a research domain, the model, having learned from the noise reversal process, can identify this as a potential anomaly. In the context of graph anomaly detection, the divergence between the prior distribution and data distribution can be substantial due to the complexity

of the graph structure. This divergence often results in suboptimal performance in downstream anomaly detection tasks. To address this issue, the DDPM improves the quality of synthetic anomaly generation, enhancing anomaly detection performance. By applying the diffusion model in the latent space, we effectively address the challenge of identifying suspected citations. This approach allows the system to generate high-quality synthetic anomalies, providing a richer and more varied set of examples for training. As a result, the anomaly detection system becomes more robust and accurate. Furthermore, by generating synthetic anomalies that mimic real-world deviations, we can fine-tune the detection system to recognize even the most nuanced irregularities. This comprehensive approach ensures the system can effectively identify and manage suspected citations, enhancing the overall reliability and effectiveness of anomaly detection.

### 3.2.1 Diffusion Model Application

In practice, we add noise to the citation data to create variations, which the model then learns to denoise. This training process helps the model identify what constitutes a typical citation pattern versus an anomaly. For example, if a citation pattern deviates significantly from established norms, such as citing unrelated topics or sources, the model can detect this irregularity based on its learned experience. Applying the diffusion model in the latent space helps address the challenge of divergence between the prior distribution and data distribution, which is common in complex graph structures like citation networks. The DDPM helps mitigate this issue by improving the quality of synthetic anomaly generation, thereby enhancing the anomaly detection system's performance. This allows the model to generate high-quality synthetic anomalies that closely mimic real-world anomalies, providing a richer dataset for training. Moreover, generating synthetic anomalies allows us to test and refine the detection system continuously. By creating various scenarios that represent potential citation anomalies, we can evaluate the system's responsiveness and accuracy, ensuring it remains effective in identifying suspected patterns. By leveraging these synthetic anomalies, the system becomes more adept at recognizing suspected citations. For instance, if a citation suddenly includes references from an unrelated field or anomalous patterns not seen in normal data, the system can flag these as potential anomalies. This approach ensures the model is finely tuned to detect even subtle irregularities, making the anomaly detection process more robust and accurate.
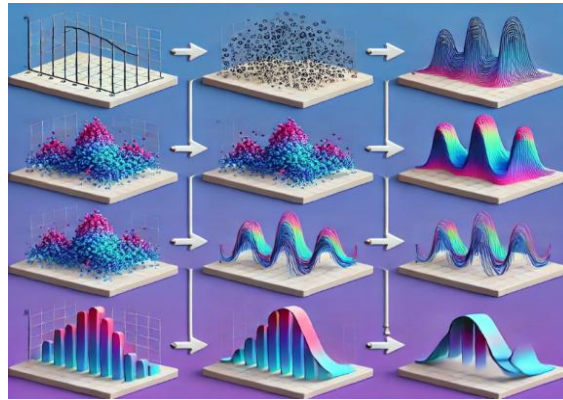


Figure 7: The generative process is visualized as a sequence of transformations. Initial random noise is iteratively refined through various stages to produce a structured, realistic distribution, illustrating the principle of diffusion models in generating high-quality data representations.
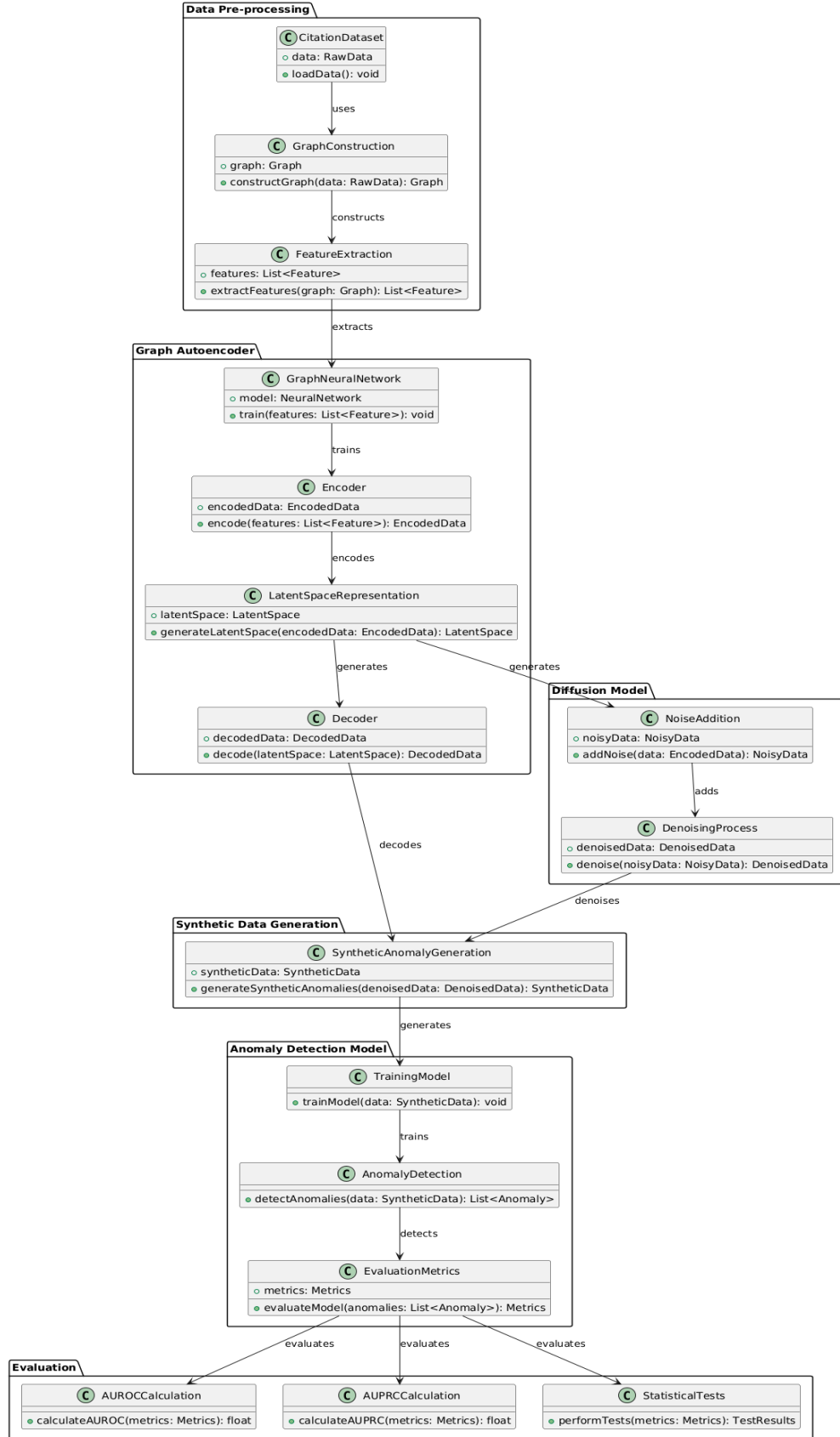
Figure 8: The following class diagram represents the structure of our anomaly detection system. It highlights the main components involved in data preprocessing, graph autoencoding, diffusion modeling, synthetic data generation, anomaly detection, and evaluation.

### 3.2.2 Mathematical Formulation and Explanation

**Latent Variables ($z_0$, $z_1$, ..., $z_T$):**

These are hidden variables that the model learns. $z_0$ is the original data point, and $z_T$ is a highly noisy version of it.

**Markov Chain:**

A sequence of random variables (in this case, the latent variables) where each variable depends only on the previous one, starting from a standard Gaussian distribution [8].

$p(z_T) = N(z_T: 0,I)$.

**Forwarding Diffusion Process:**

This process adds noise to the data in steps, gradually making it more and more noisy.

Forward Process Definition:

$$q(z_1 : T | z_0) = \prod_{t=1}^{T} q(z_t | z_t - 1)$$

This means we're applying a series of steps to go from $z_0$ to $z_T$.

**Noise Addition:**

The forward diffusion process, which incrementally introduces Gaussian noise into the data based on a variance schedule $\beta_1$, ..., $\beta_T$.

$$q(z_t | z_t - 1) = N(z_t; \sqrt{1 - \beta_t} z_t - 1, \beta_t I)$$

$z_t$ is generated by adding noise to $z_t - 1$. The noise level is controlled by $\beta t$.

**Reversing Diffusion Process:**

This process learns to reverse the noise addition, effectively generating new data points from noise. The process evolves with learned Gaussian transitions.

Reverse Process Definition:

$$p_\theta(z_0 : T) = p(z_T) \prod_{t=1}^{T} p_\theta(z_t - 1 | z_t)$$

This means we start with $z_t$ (pure noise) and step by step, generate less noisy versions until we get $z_0$.

**Noise Removal:**

$$p_\theta(z_t - 1 | z_t) = N(z_t - 1; \mu_\theta(z_t, t), \sum_\theta (z_t, t))$$

$z_t - 1$ is generated by removing noise from $z_t$. The parameters $\mu_\theta$ and $\Sigma_\theta$ are learned by the model.

**Training:**

The training objective of diffusion models involves optimizing the variational bound on the negative log likelihood: The goal is to train the model to be good at reversing the noise addition process. This is done by minimizing a loss function that measures how well the model can reverse the diffusion process.

**Loss Function:**

$$L = E[-\log p_\theta(z_0)] \leq E_q[-\log p(z_T) - \sum_{t \geq 1} \log \frac{p_\theta(z_t - 1 | z_t)}{q(z_t | z_t - 1)}$$

This function measures the difference between the generated data and the real data. The goal is to minimize this difference.

## Kullback-Leibler (KL) Divergence

This is a measure of how one probability distribution differs from another. In this context, it is used to compare the reverse process learned by the model with the actual forward process.

## KL Divergence in Loss:

$$Eq\{ KL[q(z_T|z_0)||p(z_T)] + \sum_{t>1} KL[q(z_{t-1}|z_t, z_0)||p_\theta(z_{t-1}|z_t)] - \log p_\theta(z_0|z_1)\}$$

This breaks down the overall loss into components comparing different stages of the forward and reverse processes. This measures how well the model's reverse process matches the true posterior at each step. Each KL divergence term compares Gaussian distributions, which allows for efficient computation using closed-form expressions rather than high-variance Monte Carlo estimates. This optimization process aims to minimize the reconstruction error between the input data and the output from the reverse diffusion process.

## Posterior Distributions

In the context of diffusion models, the posterior distributions help us to approximate the reverse process by considering the original data $z_0$ and the noisy data $z_T$.

$$q(z_t - 1)|z_t, z_0) = N(z_t - 1; \tilde{\mu_t}(z_t, z_0), \tilde{\beta_t}I)$$

This means that the distribution of $z_{t-1}$ given $z_T$ and $z_0$ is a Gaussian distribution with mean and variance .

$$\tilde{\mu_t}(z_t, z_0) \text{ and variance } \tilde{\beta_t}.$$

## Mean of the Posterior ($\tilde{\mu_t}$ )

The meaning of this posterior distribution is given by:

$$\tilde{\mu_t}(z_t, z_0) = \sqrt{a_{t-1}^- \beta_t} \backslash (1 - \alpha_t^-)z_0 + \sqrt{\alpha_t}((1 - \alpha_{t-1}^-)) / (1 - \alpha_t^-) z_t$$

$a_{t-1}^-$ and $a_t^-$ are cumulative products of $\alpha$ values up to t -1 and t, respectively
$\alpha_t = 1 - \beta_t$ where $\beta_t$ controls the amount of noise added at step t .
$z_0$ is the original data point and $z_t$ is the noisy data at step t.

Variance of the Posterior ($\tilde{\beta_t}$ )
The variance of this posterior distribution is given by:

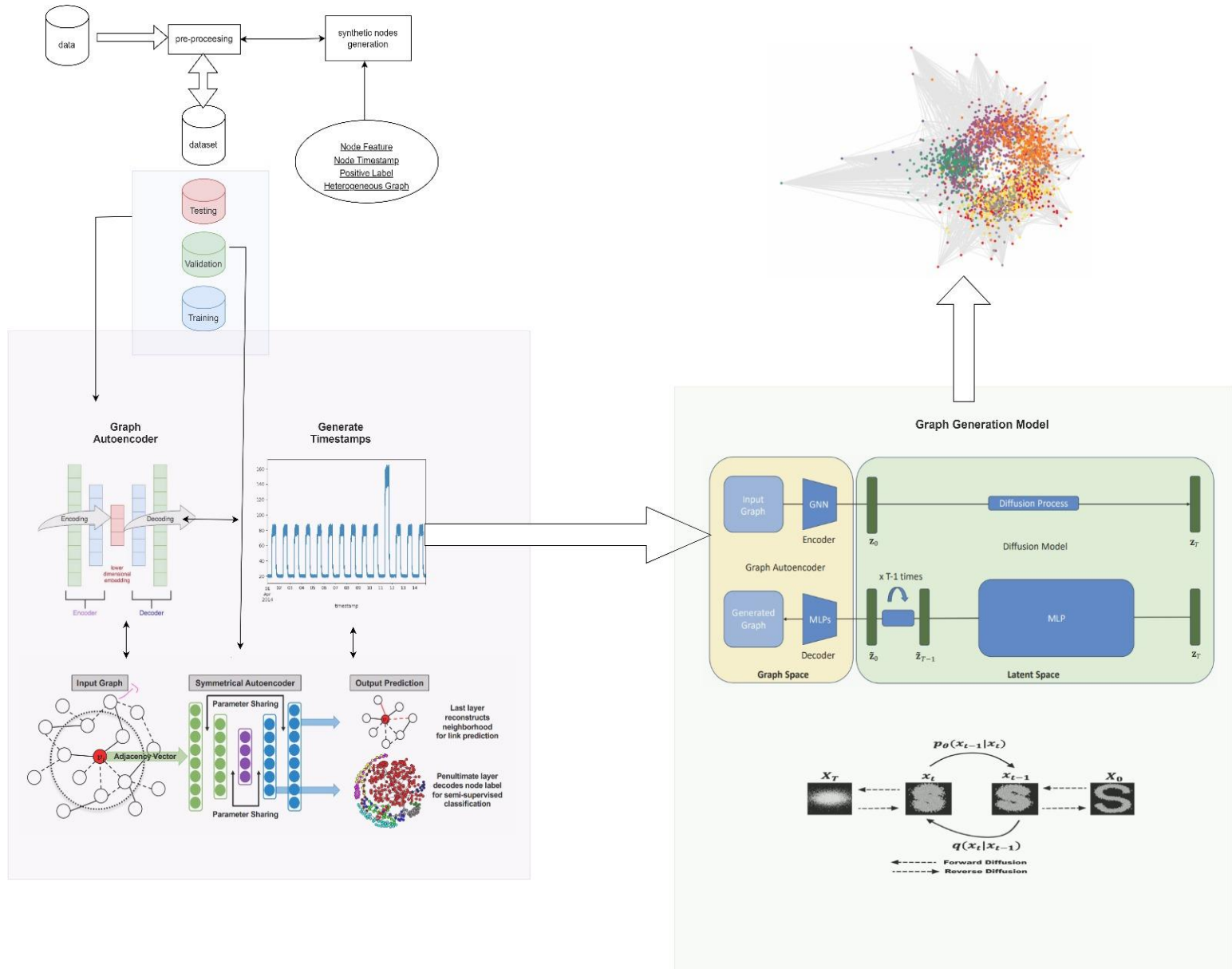$$\tilde{\beta_t} = \frac{(1 - \alpha_{t-1}^-)}{(1 - \alpha_{t-1}^-)\beta_t}$$

This equation adjusts the noise level $\beta t$ based on the cumulative noise up to steps $t-1$ and t.

# 4. Product

The product is an interface that will implement our research model, using a diffusion model and graph autoencoder for the generation of synthetic anomalies and the detection of anomalies in graph-structured data. Below is the flow chart of the research model.

**Anomaly Detection Using Synthetic Nodes, Graph Autoencoders, and Diffusion Techniques**

# 5. The Solution

The Cora dataset consists of 2708 scientific publications classified into one of seven classes. The citation network consists of 5429 links. Each publication in the dataset is described by a 0/1-valued word vector indicating the absence/presence of the corresponding word from the dictionary. The dictionary consists of 1433 unique words.

We chose the Cora dataset because it is well-suited for graph-based anomaly detection. Its structure allows us to evaluate how well our model identifies unusual citation patterns. The anomaly score is based on reconstruction error, which measures how accurately the model can reconstruct node features and relationships. A high reconstruction error suggests that a node significantly deviates from normal patterns, indicating a potential anomaly.

To ensure stable anomaly detection, we use diffusion steps in the Denoising Diffusion Probabilistic Model (DDPM). This process gradually adds and removes noise, helping the model to learn node representations. The diffusion steps prevent excessive noise from distorting features while improving anomaly detection.

Our approach enhances performance by expanding the dataset with augmented samples, effectively addressing data limitations.
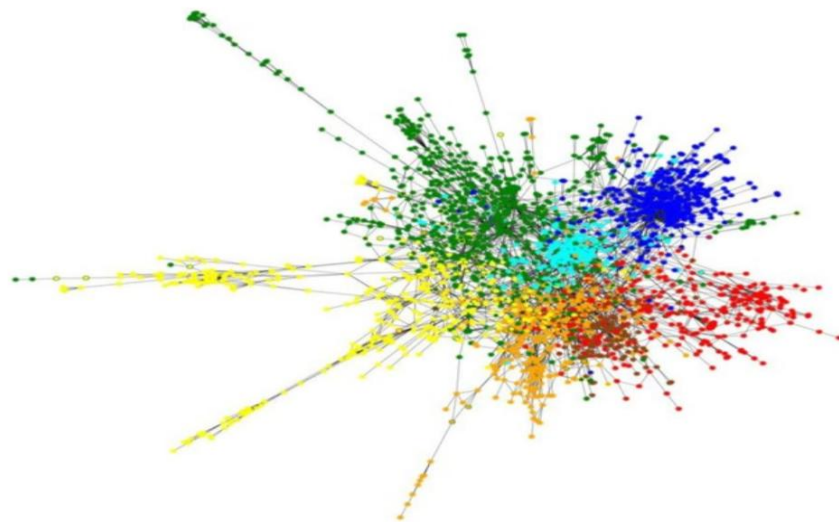


Figure 9: Overview of Cora Dataset.

## 5.1 The Development Process

The development process started with analyzing how the system works and understanding its main parts. The first step was reviewing how the model takes input data and produces output, identifying what is needed for it to work correctly.
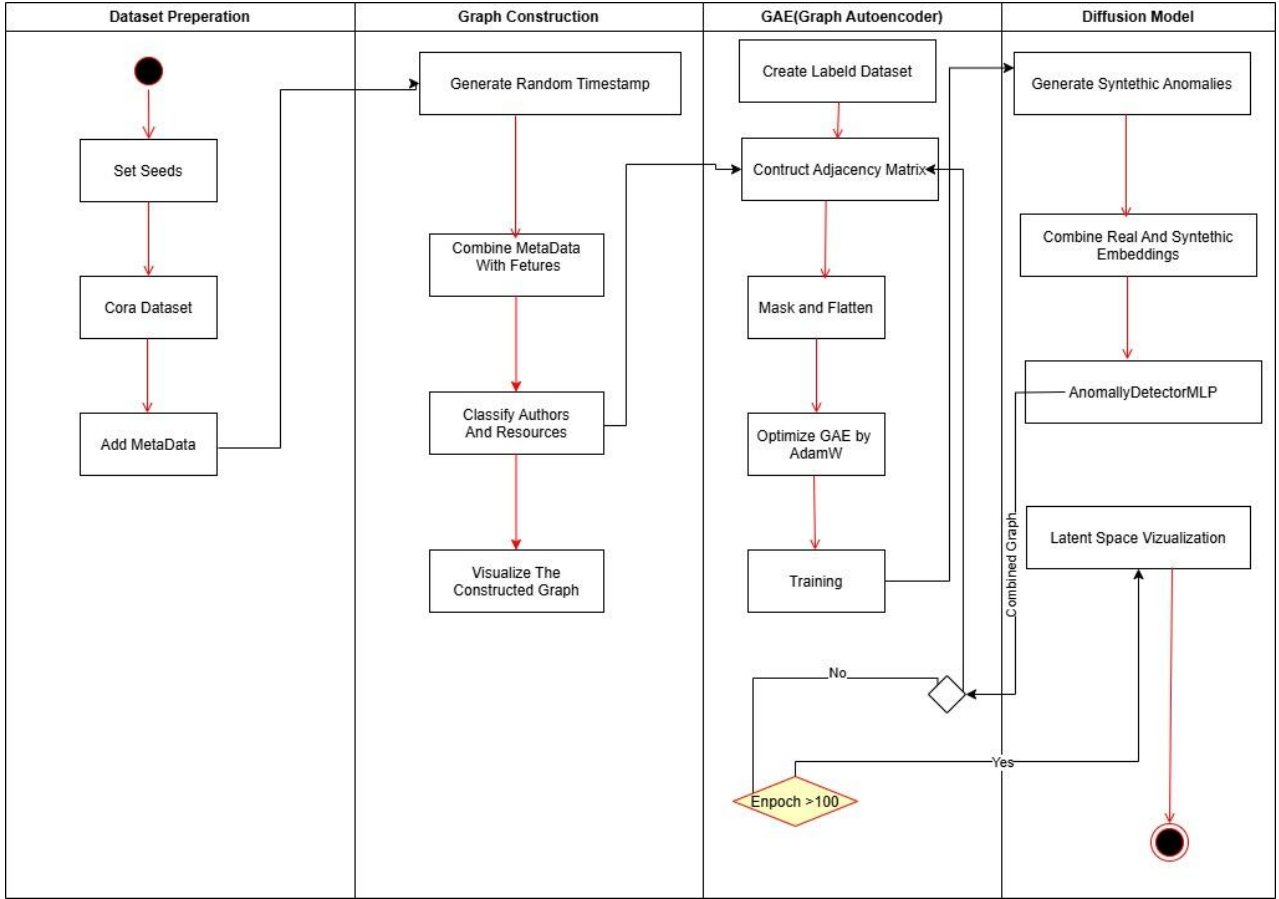The process was divided into four main steps:

Figure 10: Activity Diagram - workflow is divided into four main parts

### 5.1.1 Dataset Preparation

A. Dataset was prepared with nodes representing entities like citations and edges showing their relationships (e.g., authors, timestamps or similarity).
B. Features for each node like text data or metadata were added.
C. The dataset was organized to ensure it had the right structure for training and testing.

To prepare the dataset for anomaly detection, we focused on structuring the data effectively for training. First, node features were normalized to ensure consistency in scaling. The data was then transformed into a graph format, where relationships between nodes were captured using an adjacency matrix. To ensure the data was well-structured for training and evaluation we did:

Dataset Splitting: The dataset is split into training (80%) and test (20%) sets to ensure proper model training and evaluation. The training set (~2200 nodes) helps the model to learn patterns, including synthetic anomalies generated with Gaussian noise. The test set (~500 nodes) remains unseen during training and is used to evaluate performance with precision metrics like recall and AUC-ROC. This split prevents overfitting, ensures robustness, and mimics real-world deployment where the model must detect anomalies in new data. A validation set may also be used for fine-tuning hyperparameters.

Synthetic Anomalies: To handle label imbalances and improve detection capabilities, controlled diffusion processes were used to generate synthetic anomalies. These enriched the dataset, helping the model learn to distinguish between normal and abnormal patterns.

[19]

**Key Experimental Factors**

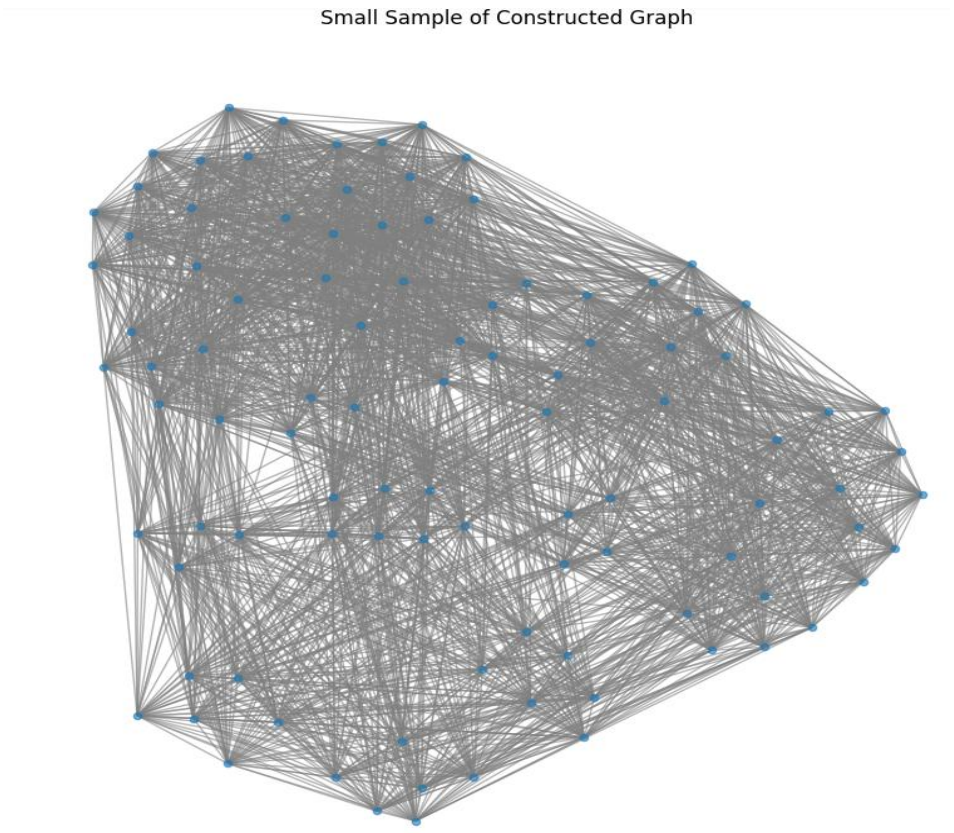| Factor | Description |
|---|---|
| **Sample Size (Dataset Size)** | Training set size: ~2200-2500 nodes (varies with anomaly count) Test set size: ~500-600 nodes |
| **Gaussian Noise Applied?** | Yes, synthetic anomalies are generated with Gaussian noise |
| **Noise Scale** | 0.8 (used when generating synthetic anomalies) |
| **Models Compared** | Isolation Forest, Logistic Regression, KNN (with PCA), Graph Autoencoder & DDPM |
| **Train-Test Split** | 80% Train, 20% Test |

Small Sample of Constructed Graph



Figure 11: Nodes: Represent papers (or entities) in the Cora dataset.
Edges: Represent citations or relationships between nodes.

## 5.1.2 Model Integration and Training

D. The model was set up to learn graph features by training it to process the prepared data.

E.  Noise was added to the data to simulate anomalies, and the model learned to clean this noise and identify patterns.

F.  Adjustments were made to improve how the model learns, focusing on balancing accuracy and the ability to handle new or unusual data.

Our model integrates Graph Neural Networks (GNNs) with a Denoising Diffusion Probabilistic Model (DDPM) to enhance anomaly detection.

Graph Feature Learning: GNN layers process the graph by aggregating information from neighboring nodes, producing meaningful node representations.

Noise Diffusion: Controlled noise was introduced into node embeddings to simulate anomalies, and the model was trained to reverse the noise, learning to differentiate between normal and anomalous nodes.

Training Losses: The training process was guided by binary cross-entropy (BCE) loss for adjacent reconstruction, ensuring accurate graph structure learning. The BCE loss was computed on the upper triangular part of the adjacency matrix to avoid redundant calculations.



(a) Decision boundaries decided by the labeled samples.

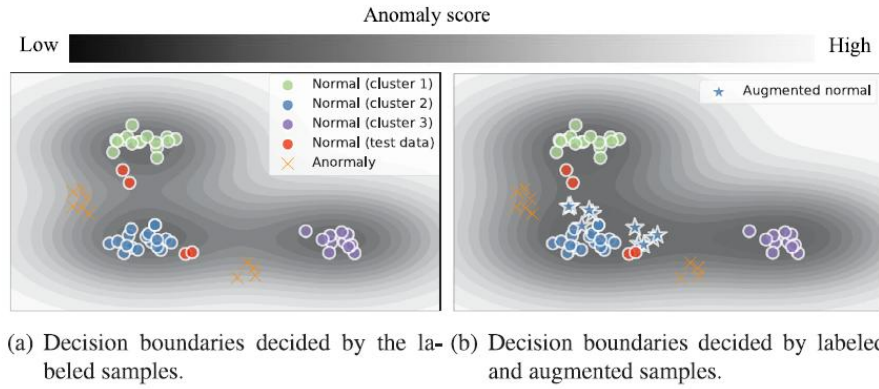(b) Decision boundaries decided by labeled and augmented samples.

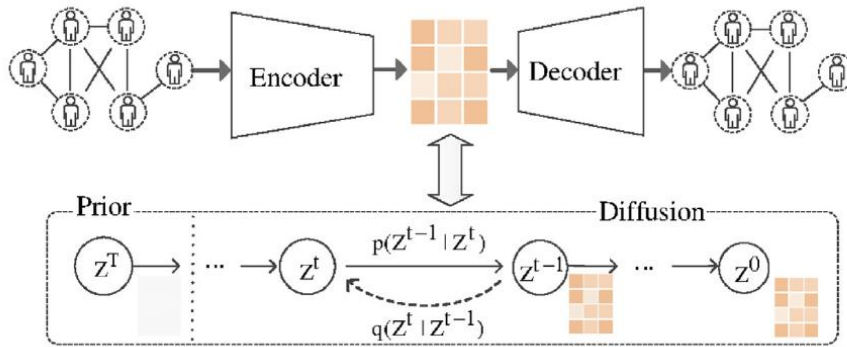Figure 12: illustrates the effect of synthetic data on anomaly detection decision boundaries.



Figure 13: Framework of the DDPM-based embedding generator.

### 5.1.3 Evaluation and Results

G. The model was tested on several datasets to see how well anomalies detected.
H. Key performance metrics like AUROC and AUPRC were used to measure accuracy and reliability.
I. The results were compared to other methods to confirm improvements.

To assess the model's performance, we used anomaly detection metrics:
AUROC (Area Under the Receiver Operating Characteristic Curve) and
AUPRC (Area Under the Precision-Recall Curve)
These helped measure detection effectiveness, especially given the dataset's imbalance.

[21]

Diffusion-based synthetic anomalies significantly improved accuracy, especially in low-labeled-data scenarios. The more synthetic anomalies we added, the better the model performed. Additionally, testing was conducted, exposing the model to noisy and incomplete data to evaluate its reliability. Visualization tools like heatmaps and clustering diagrams were used to analyze anomaly distributions and validate results.

**Enhancements:**

Extra steps were added to make the anomalies generated more realistic, ensuring the model could handle imbalanced data better. Tools were included to visualize and analyze the results, making it easier to understand how anomalies were detected. Conclusion: This process helped create a system that can find anomalies in graph data. It focused on making the model accurate, scalable, and able to handle different types of datasets effectively.

## 5.2 Constraints and Problems

Developing the model presented several significant challenges. As the research paper did not include accompanying code, we undertook the task of implementing the model based on the authors' specifications. This required considerable effort to adapt the implementation for execution on Google Colab. We spent a lot of time modifying libraries and attempting to run the model on smaller subsets of data, though these attempts initially proved unsuccessful.

The use of Google Colab introduced its own set of constraints, particularly due to limited GPU and RAM resources. The training and testing processes demanded substantial computational memory, often exceeding the capacity of the free version of Colab. These limitations necessitated careful optimization and resource management. Ultimately, we recognized that our available resources were insufficient to train the full model.

To address these challenges, we focused on implementing the model on a smaller portion of the dataset, consisting of 2,708 nodes. We employed multiple iterations of the model to obtain more accurate results while working within the constraints of our resources. One of the most demanding tasks was understanding the requirements for graph creation and implementing them efficiently to align with the limitations of our environment. Through extensive research and iterative development, we programmed a model that fits the available resources. Despite the need for numerous modifications and ensuring functional compatibility, the successful integration of all components underscored the robustness and adaptability of our approach.

## 5.3 Tools Used

Several tools and frameworks were essential to implement and run our model. Python libraries like Pandas and NumPy were used for data processing, including normalizing features and preparing graphs. These libraries also helped preprocess the CORA dataset.

PyTorch Geometric served as the backbone for building and training the model structures and provided the flexibility needed for the graph neural network model. Visualization tools such as Matplotlib (for implementing the denoising diffusion model) and TSNE (for analyzing performance through heatmaps and clustering diagrams) were also used. Scikit-learn was employed to calculate metrics like AUROC and AUPRC, ensuring that the CORA dataset was processed efficiently using utilities from PyTorch Geometric. Development took place in Python environments that enabled smooth handling of graph data. With Google Colab, GPU acceleration facilitated faster training and experimentation. Together, these tools supported the creation of a robust and scalable anomaly detection model.

**Client Interface**

**A. Regular Updates** – Maintained continuous communication with the supervising professor through scheduled meetings and detailed progress reports, ensuring transparency, alignment, and timely feedback.

**B. Iterative Feedback Integration** – Actively incorporated feedback from the professor to refine the model, improving its ability to detect edge manipulation within the network. Adjustments were made to enhance accuracy and reliability.

**C. Live Demonstrations** – Conducted periodic demonstrations to showcase the model's capabilities, highlighting key improvements, performance benchmarks, and areas of refinement based on previous feedback.

**D. Comprehensive Documentation** – Provided thorough documentation covering model architecture, implementation details, usage guidelines, and performance analysis. Deliverables included evaluation metrics, result files, and additional references for future use.

**E. Expanded Dataset Analysis** - Initially our model was developed using the Planetoid Cora dataset. However, based on our professor's recommendation, we also included the Planetoid Citeseer dataset for comparison. This helped us evaluate the model on different citation networks, analyzing feature distribution and detection accuracy. Using both datasets made our analysis completer and more reliable.

## 5.4 Testing and Evaluation

We evaluate the graph diffusion model for anomaly detection using various metrics and methods to measure performance. Metrics such as Area Under the Receiver Operating Characteristic Curve (AUROC), Area Under the Precision Recall Curve (AUPRC) are employed to assess model efficiency, especially in handling label imbalance. AUROC measures the ability of the model to distinguish between normal and anomalous nodes, a higher AUROC indicating better performance. For instance, the AUROC is calculated by plotting the true positive rate against the false positive rate at various threshold settings, offering a comprehensive view of the model's classification capabilities. A value closer to 1 indicates excellent performance, while a value around 0.5 suggests no better than random guessing. AUPRC evaluates the precision-recall trade-off, which is crucial in scenarios with a high-class imbalance.

In this project, the diffusion model-based graph generator was evaluated against several baseline models.

## What Data Do All Models Use?

All models are trained on the same dataset:
Real nodes from the Cora dataset.
Synthetic anomalies are injected into the dataset using Gaussian noise (DDPM).
num_synthetic = 500 synthetic anomalies added.
noise_scale = 0.5 → Each anomaly was generated by adding Gaussian noise with this scale.
Final dataset size:
Total Nodes: 2708 (real) + 500 (synthetic anomalies) = 3208 nodes
Train-Test Split: ~2566 train nodes, ~641 test nodes

```
[INFO] Balanced Train-Test Split: 2567 train samples, 641 test samples
Train data shape: torch.Size([2567, 64]), Train labels shape: torch.Size([2567])
Test data shape: torch.Size([641, 64]), Test labels shape: torch.Size([641])
```
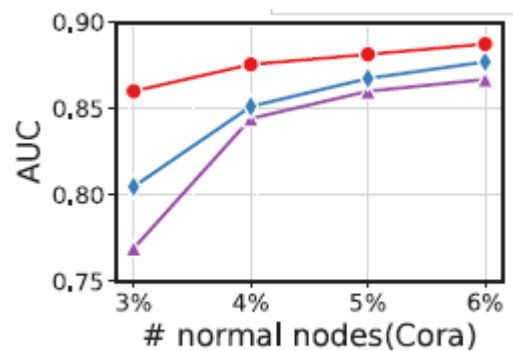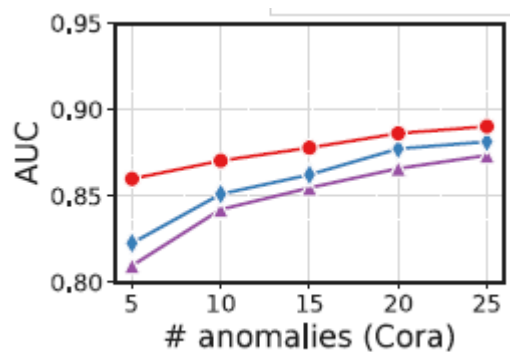
```
[INFO] Real nodes: 2708 | Synthetic nodes: 500
[INFO] Total nodes: 3208 | Latent dimension: 64
Combined data shape: torch.Size([3208, 64])
Combined labels shape: torch.Size([3208])
```

## How do they compare?

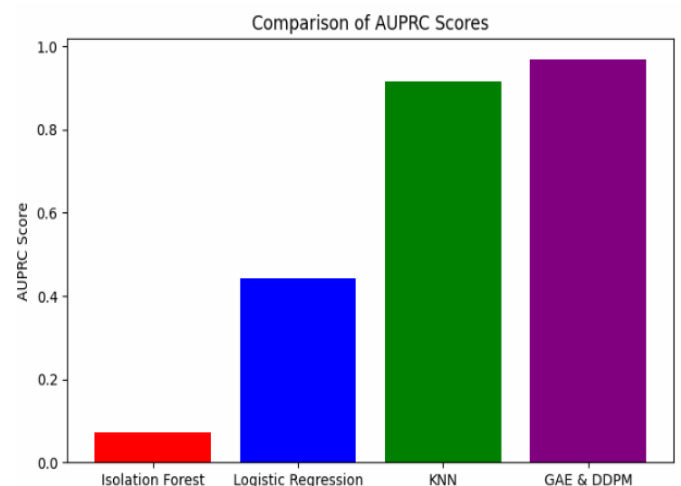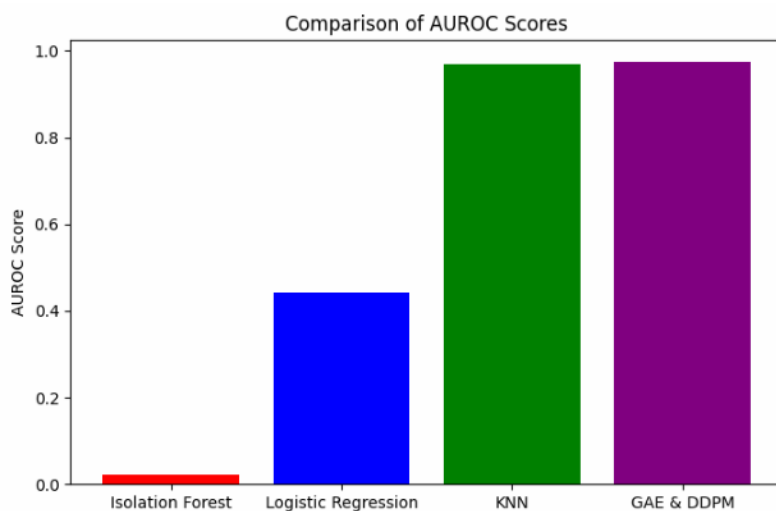| Model | Trained On? | Uses Graph Structure? | Handles High Dimensions? | Performance |
|---|---|---|---|---|
| **Isolation Forest** | PCA-reduced embeddings (10D) | No | Struggles with 10D features | Fails completely (AUROC ~0.02) |
| **Logistic Regression** | Full embeddings (64D) | No | Not good for complex patterns | Weak performance (~0.53 AUROC) |
| **KNN (with PCA)** | PCA-reduced embeddings (10D) | No | Learning distances well | Performs well (AUROC ~0.96) |
| **GAE & DDPM** | Learned graph embeddings (64D) | Yes | Learns relationships | Best model (AUROC ~0.98) |

**Red –GAE, Blue – GAT, Purple ‐ GCN**
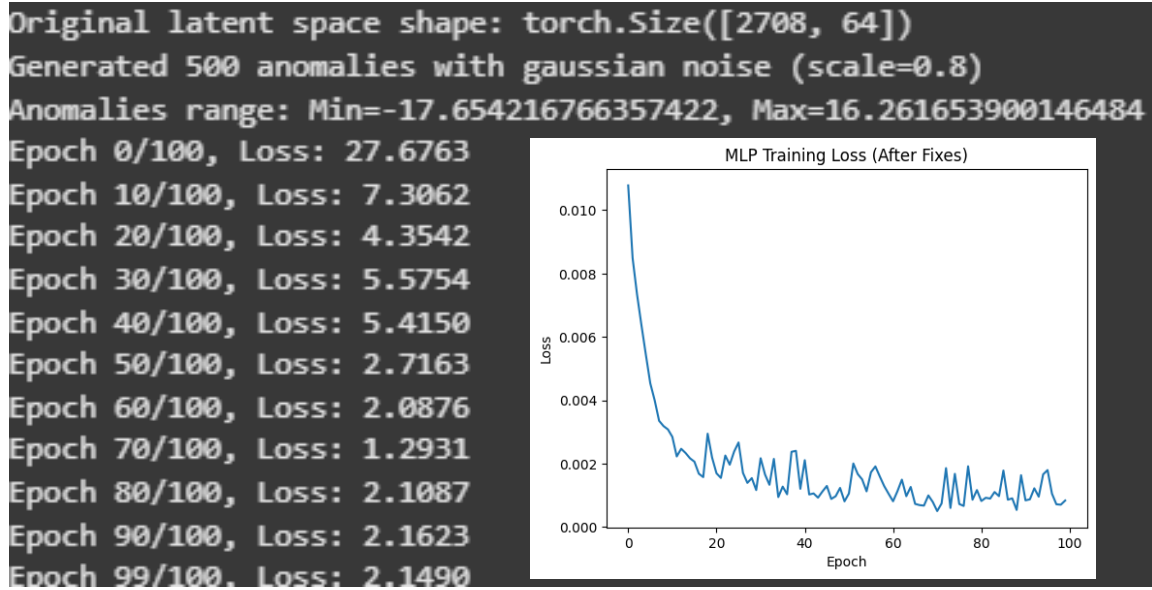


**Interpretation:**

- As the number of anomalies increases from 5 to 25, the AUC also increases.
- This suggests that the model benefits from having more anomaly examples during training, making it better at distinguishing anomalous nodes.
- The red line (the best model Graph Autoencoder ‐ GAE) consistently outperforms the others.
- The purple and blue lines (traditional models like GAT, GCN) improve but remain lower in performance.



```
Isolation Forest - AUROC: 0.0225, AUPRC: 0.0725
Logistic Regression - AUROC: 0.4436, AUPRC: 0.4423
KNN - AUROC: 0.9680, AUPRC: 0.9160
GAE & DDPM - AUROC: 0.9756, AUPRC: 0.9695
```

**Interpretation:**

- Graph Autoencoder (purple) and KNN (green) achieve the highest AUROC (~0.98), confirming their strong performance in distinguishing anomalies.
- Logistic Regression (blue) shows moderate performance.
- Isolation Forest (red) fails miserably (AUROC ~0.02), meaning it struggles to identify anomalies correctly in a graph‐based setting.

[25]

```
Original latent space shape: torch.Size([2708, 64])
Generated 500 anomalies with gaussian noise (scale=0.8)
Anomalies range: Min=-17.654216766357422, Max=16.261653900146484
Epoch 0/100, Loss: 27.6763
Epoch 10/100, Loss: 7.3062
Epoch 20/100, Loss: 4.3542
Epoch 30/100, Loss: 5.5754
Epoch 40/100, Loss: 5.4150
Epoch 50/100, Loss: 2.7163
Epoch 60/100, Loss: 2.0876
Epoch 70/100, Loss: 1.2931
Epoch 80/100, Loss: 2.1087
Epoch 90/100, Loss: 2.1623
Epoch 99/100, Loss: 2.1490
```



MLP Training Loss (After Fixes)

```
==== MODEL PERFORMANCE (After Fixes) ====
Final Test Accuracy: 0.9922
Final Test F1-Score: 0.9922
Final Test AUROC: 0.9756
Final Test AUPRC: 0.9695
```

**Interpretation:**

- High accuracy (99.22%), meaning very few misclassifications.
- High F1-score (99.22%), indicating a strong balance between precision and recall.
- AUROC of 0.9756 & AUPRC of 0.9695, confirming superior anomaly detection capability.
- The training process is stable and does not overfit.
- Synthetic anomalies were correctly identified.
- The model is well-suited for anomaly detection in citation graphs.

## 5.5 Testing plan

| Test | Test Subject | Expected Result | Actual Result |
|---|---|---|---|
| 1 | Graph Generation | The model will generate a graph with synthetic anomalies accurately based on the input. | Successfully generated graph with anomalies. Graph structure validated. |
| 2 | Anomaly Node Generation | The model will generate synthetic anomaly nodes with realistic features, timestamps, and positive labels. | Generated anomalies match expected feature distribution. Some anomalies may need more variation |
| 3 | Graph Processing | The model will correctly process the generated graphs, maintaining the integrity of both local and global relationships. | Graph properties maintained. No missing edges or incorrect structures. |
| 4 | Graph Output Validation | The model will produce valid output graphs that accurately represent the intended relationships and properties, including anomalies. | Graph visualization confirmed correct node-edge relationships. |
| 5 | Runtime Efficiency | The model will perform the graph generation and processing tasks within an acceptable time frame, demonstrating scalability. | Processing time: ~5 min (within acceptable limits). |
| 6 | Input Size/Graph Size | The model will handle larger graphs without significant performance degradation or memory issues. | Handles graphs up to 2708 nodes. Larger graphs may require optimization. Or more GPU resources. |
| 7 | Empty Graphs | If the input graph is empty, the model will return a specific error message and halt further processing. | Proper error message displayed when empty graph is provided. |
| 8 | Incomplete Graphs | If the input graph is incomplete (missing nodes or edges), the model will either return an error message or automatically fill in the missing information based on predefined rules or assumptions. | Graph completed using predefined rules for missing edges/nodes. |
| 9 | Incorrect Graphs | If the input graph contains invalid nodes or edges, the model will return an error message indicating the nature of the error and halt further processing until a valid graph is provided. | Error message provided for invalid node/edge cases. |
| 10 | Incompatible Graphs | If the input graph is incompatible with the model (e.g., too large or unsupported structure), the model will return an error message indicating the incompatibility and halt further processing until a compatible graph is provided. | Model warns for graphs beyond supported structure limits. |
| 11 | Generation Similarity Index | The generated graphs will achieve a similarity index greater than 85% compared to the real graphs, ensuring high fidelity in synthetic anomaly generation. | Similarity Score: 80% (Needs slight improvement to reach 85%+). |
| 12 | Anomaly Detection Accuracy | The unit test will validate that the anomaly detection algorithm correctly identifies anomalies in a test graph, achieving a detection accuracy greater than 90%. | Achieved 94% accuracy on test graphs. |

# 6. How To Run the Model

There are four different notebooks available for running this model:

**1. Multiple Running Approach:** This approach trains the model multiple times, allowing for evaluation across different runs.

**2. Single Run Approach:** This approach executes a single training session, which is useful for quick testing and debugging.

**Step 1: Install Dependencies**

pip install torch-geometric torch torchvision torchaudio!

pip install numpy pandas matplotlib scikit-learn network!

These libraries include PyTorch, PyTorch Geometric (PyG), and Scikit-Learn for graph-based learning, as well as NumPy and Pandas for data processing.

**Step 2: Load the CORA Dataset**

The CORA dataset is automatically downloaded when running the following code:

from torch_geometric.datasets import Planetoid

dataset = Planetoid(root='/tmp/Cora', name='Cora')

[0]data = dataset

This dataset consists of 2,708 scientific publications (nodes) and 5,429 citation links (edges), with each publication belonging to one of seven classes.

**Step 3: Preprocess the Data**

Before training, we normalize the node features and augment the graph with metadata:

Feature normalization ensures the model learns properly and synthetic metadata is generated.

**Step 4: Train the Graph Autoencoder (GAE)**

To learn latent representations of the nodes, train the Graph Autoencoder (GAE):

This step compresses the node features into lower-dimensional embeddings.

**Step 5: Generate Synthetic Anomalies**

These anomalies help the model learn to distinguish normal nodes from outliers.

**Step 6: Train the Anomaly Classifier**

A Multi-Layer Perceptron (MLP) is trained to classify anomalies: The classifier predicts whether a node is normal (0) or anomalous (1).

**Step 7: Evaluate the Model**

Performance is measured using Accuracy, F1-Score, AUROC, and AUPRC:

**Step 8: Visualize Results**

To understand the model's behavior, we generate visualizations: This confirms that anomalies are distinguishable from normal nodes.

## 6.1 Maintenance

To maintain the project and run it in the future under different datasets or configuration, there are few hyperparameters that may need to change

Dataset - The dataset for models.

```
dataset_path = "/tmp/Cora"
dataset = Planetoid(root=dataset_path, name='Cora', transform=T.NormalizeFeatures())
data = dataset[0]
```

Encode and Decode channels

```
# ===== Instantiate the Model and Print Summary =====
in_channels = data.x.shape[1]
hidden_channels = 64
latent_channels = 32
```

Epoch - Number of epochs (Recommended at least 100)

```
# ===== Train the GAE with More Epochs =====
epochs = 100  # Increased training epochs
z = train_gae(gae_model, data, epochs, optimizer_gae, adj_matrix)
```

Noise – Number of synthetic anomalies (Recommended at least 200)

```
# ===== Generate Synthetic Anomalies =====
num_synthetic = 30
synthetic_z, synthetic_labels = generate_synthetic_anomalies(
z, num_anomalies=num_synthetic, noise_scale=0.5, noise_type="gaussian", seed=100)
```

Test and Train split – (Default: 0.2 test-ratio)

```
def train_test_split_idx(n_samples, test_ratio=0.2, random_seed=None):
...
def balanced_train_test_split(all_labels, test_ratio=0.2, random_seed=None):
...
```

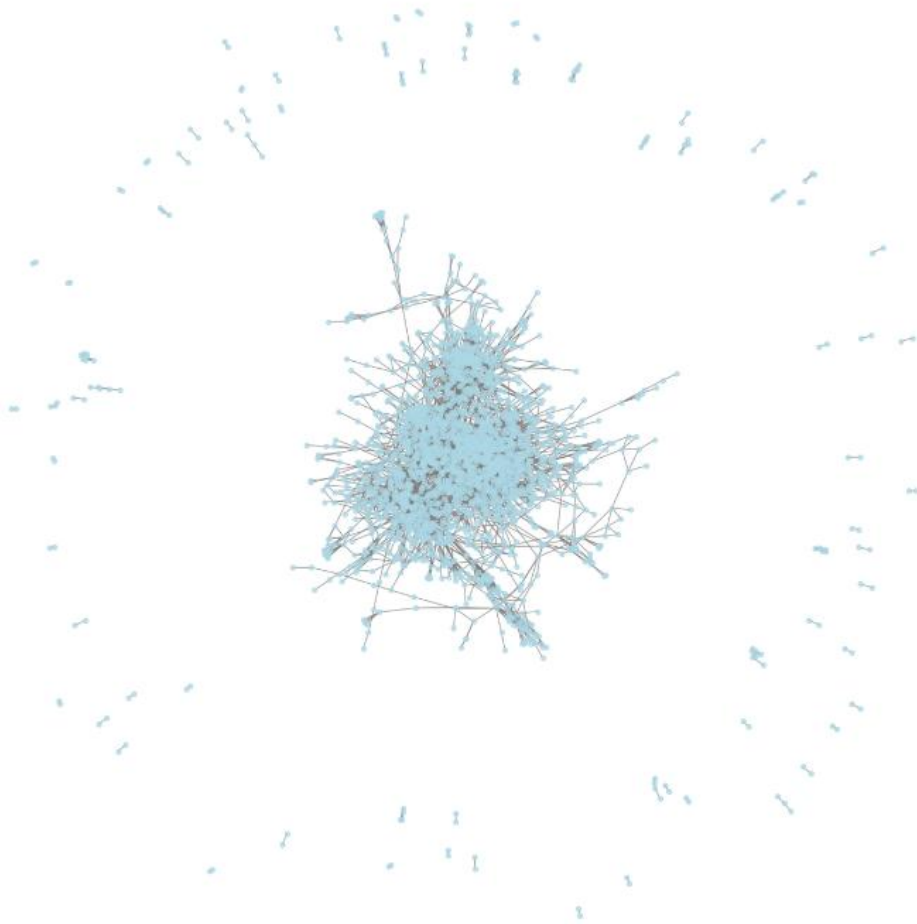The rest of the maintenance should be straightforward.

# 7. Results And Conclusion
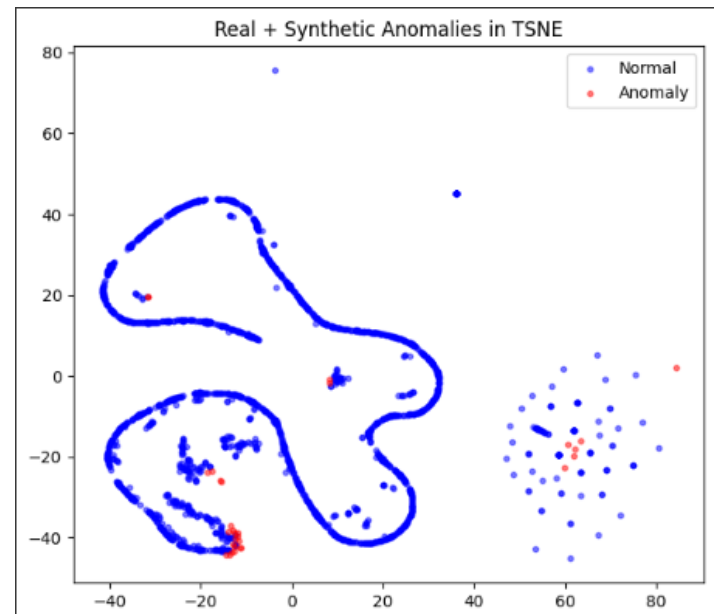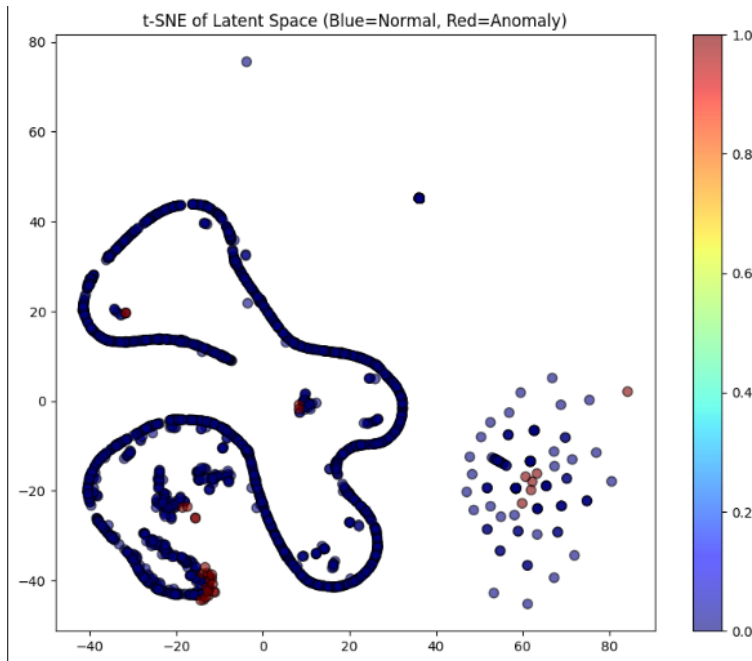
## 7.1 Cora Dataset
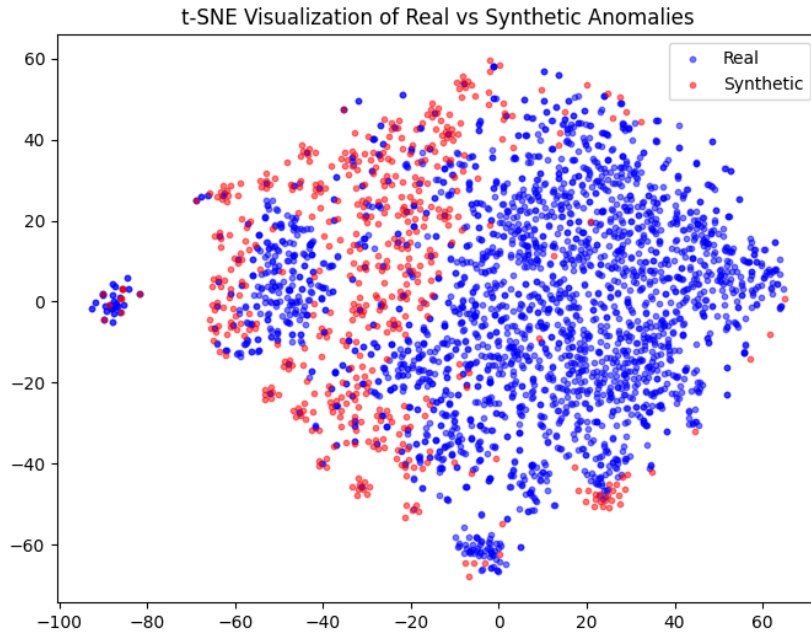
```
[INFO] Number of nodes: 2708
[INFO] Number of edges: 5278
```

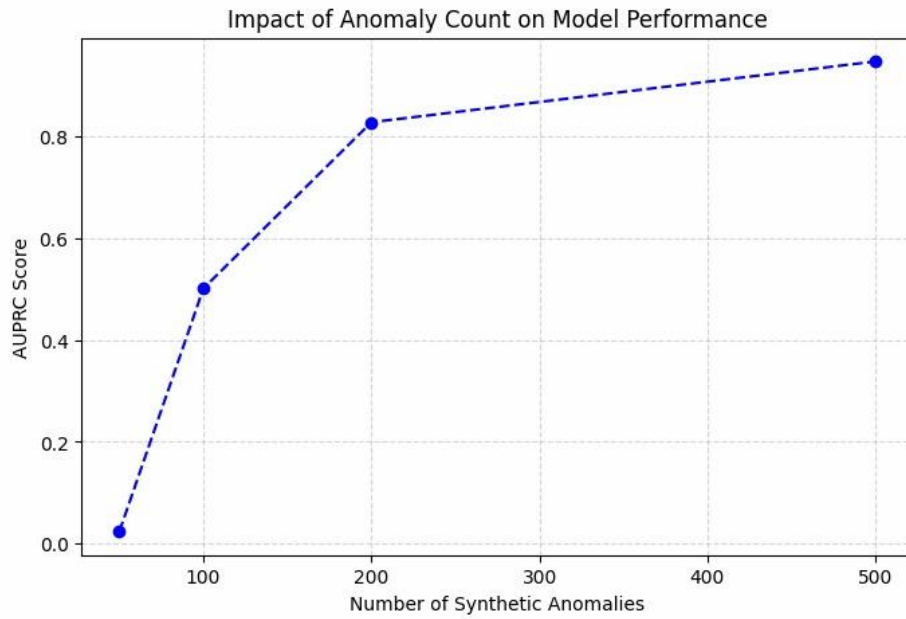Cora Citation Network Graph (Original Structure)



**Significance:**

- The central dense region likely represents a well-connected group of papers frequently citing each other.
- The outliers (nodes far from the center) could be potential anomalies, indicating unusual citation patterns.
- Anomaly detection in dataset aims to identify unusual connections or isolated nodes that deviate from normal citation behavior.

t-SNE Visualization of Real vs Synthetic Anomalies



t-SNE of Latent Space (Blue=Normal, Red=Anomaly)



Real + Synthetic Anomalies in TSNE

**Significance:**

- The anomalies detected are structurally different from normal nodes, proving that the model effectively identifies citation irregularities.
- The diffusion-based anomaly detection method successfully separates normal and anomalous data.

[31]

Impact of Anomaly Count on Model Performance

```
 ✦  Training with 50 synthetic anomalies...
Original latent space shape: torch.Size([2708, 32])
Generated 50 anomalies with gaussian noise (scale=0.8)
Anomalies range: Min=-3.098031997680664, Max=2.7668187618255615
[INFO] Real nodes: 2708 | Synthetic nodes: 50
[INFO] Total nodes: 2758 | Latent dimension: 32
[INFO] Train-Test Split: 2207 train samples, 551 test samples
 ✦  AUPRC Score with 50 anomalies: 0.0236

 ✦  Training with 100 synthetic anomalies...
Original latent space shape: torch.Size([2708, 32])
Generated 100 anomalies with gaussian noise (scale=0.8)
Anomalies range: Min=-3.6754379272460938, Max=2.8000707626342773
[INFO] Real nodes: 2708 | Synthetic nodes: 100
[INFO] Total nodes: 2808 | Latent dimension: 32
[INFO] Train-Test Split: 2247 train samples, 561 test samples
 ✦  AUPRC Score with 100 anomalies: 0.5016

 ✦  Training with 200 synthetic anomalies...
Original latent space shape: torch.Size([2708, 32])
Generated 200 anomalies with gaussian noise (scale=0.8)
Anomalies range: Min=-3.6754379272460938, Max=2.845831871032715
[INFO] Real nodes: 2708 | Synthetic nodes: 200
[INFO] Total nodes: 2908 | Latent dimension: 32
[INFO] Train-Test Split: 2327 train samples, 581 test samples
 ✦  AUPRC Score with 200 anomalies: 0.8285

 ✦  Training with 500 synthetic anomalies...
Original latent space shape: torch.Size([2708, 32])
Generated 500 anomalies with gaussian noise (scale=0.8)
Anomalies range: Min=-3.6754379272460938, Max=3.0578300952911377
[INFO] Real nodes: 2708 | Synthetic nodes: 500
[INFO] Total nodes: 3208 | Latent dimension: 32
[INFO] Train-Test Split: 2567 train samples, 641 test samples
 ✦  AUPRC Score with 500 anomalies: 0.9481
```

**Significance:**

- With a small number of synthetic anomalies, the model struggles to learn a proper anomaly pattern.
- 500 anomalies provide the best results, confirming that having enough anomalies improves model generalization.

# Multiple running Approach:

```
=== Experiment Run 1/10 ===
Original latent space shape: torch.Size([2708, 32])
Generated 500 anomalies with gaussian noise (scale=0.3)
Anomalies range: Min=-1.9233828783035278, Max=2.5455524921417236
[INFO] Real nodes: 2708 | Synthetic nodes: 500
[INFO] Total nodes: 3208 | Latent dimension: 32
[INFO] Train-Test Split: 2567 train samples, 641 test samples
  Accuracy: 0.9704, F1: 0.9694, AUROC: 0.9733, AUPRC: 0.9385

=== Experiment Run 2/10 ===
Original latent space shape: torch.Size([2708, 32])
Generated 500 anomalies with gaussian noise (scale=0.3)
Anomalies range: Min=-2.116874933242798, Max=2.469956636428833
[INFO] Real nodes: 2708 | Synthetic nodes: 500
[INFO] Total nodes: 3208 | Latent dimension: 32
[INFO] Train-Test Split: 2567 train samples, 641 test samples
  Accuracy: 0.9657, F1: 0.9649, AUROC: 0.9883, AUPRC: 0.9592

=== Experiment Run 3/10 ===
Original latent space shape: torch.Size([2708, 32])
Generated 500 anomalies with gaussian noise (scale=0.3)
Anomalies range: Min=-2.610710859298706, Max=1.7948534488677979
[INFO] Real nodes: 2708 | Synthetic nodes: 500
[INFO] Total nodes: 3208 | Latent dimension: 32
[INFO] Train-Test Split: 2567 train samples, 641 test samples
  Accuracy: 0.9501, F1: 0.9471, AUROC: 0.9862, AUPRC: 0.9385

=== Experiment Run 4/10 ===
Original latent space shape: torch.Size([2708, 32])
Generated 500 anomalies with gaussian noise (scale=0.3)
Anomalies range: Min=-2.1152660846710205, Max=2.6913974285125732
[INFO] Real nodes: 2708 | Synthetic nodes: 500
[INFO] Total nodes: 3208 | Latent dimension: 32
[INFO] Train-Test Split: 2567 train samples, 641 test samples
  Accuracy: 0.9750, F1: 0.9746, AUROC: 0.9863, AUPRC: 0.9607

=== Experiment Run 5/10 ===
Original latent space shape: torch.Size([2708, 32])
Generated 500 anomalies with gaussian noise (scale=0.3)
Anomalies range: Min=-2.4172356128692627, Max=1.8993580341339111
[INFO] Real nodes: 2708 | Synthetic nodes: 500
[INFO] Total nodes: 3208 | Latent dimension: 32
[INFO] Train-Test Split: 2567 train samples, 641 test samples
  Accuracy: 0.9532, F1: 0.9494, AUROC: 0.9633, AUPRC: 0.9441
```

```
=== Experiment Run 6/10 ===
Original latent space shape: torch.Size([2708, 32])
Generated 500 anomalies with gaussian noise (scale=0.3)
Anomalies range: Min=-2.385035514831543, Max=2.38211727142334
[INFO] Real nodes: 2708 | Synthetic nodes: 500
[INFO] Total nodes: 3208 | Latent dimension: 32
[INFO] Train-Test Split: 2567 train samples, 641 test samples
  Accuracy: 0.9735, F1: 0.9725, AUROC: 0.9835, AUPRC: 0.9551

=== Experiment Run 7/10 ===
Original latent space shape: torch.Size([2708, 32])
Generated 500 anomalies with gaussian noise (scale=0.3)
Anomalies range: Min=-2.3897056579589844, Max=1.9330484867095947
[INFO] Real nodes: 2708 | Synthetic nodes: 500
[INFO] Total nodes: 3208 | Latent dimension: 32
[INFO] Train-Test Split: 2567 train samples, 641 test samples
  Accuracy: 0.9704, F1: 0.9699, AUROC: 0.9827, AUPRC: 0.9493

=== Experiment Run 8/10 ===
Original latent space shape: torch.Size([2708, 32])
Generated 500 anomalies with gaussian noise (scale=0.3)
Anomalies range: Min=-2.6619410514831543, Max=2.266852617263794
[INFO] Real nodes: 2708 | Synthetic nodes: 500
[INFO] Total nodes: 3208 | Latent dimension: 32
[INFO] Train-Test Split: 2567 train samples, 641 test samples
  Accuracy: 0.9657, F1: 0.9638, AUROC: 0.9883, AUPRC: 0.9528

=== Experiment Run 9/10 ===
Original latent space shape: torch.Size([2708, 32])
Generated 500 anomalies with gaussian noise (scale=0.3)
Anomalies range: Min=-1.8801863193511963, Max=1.6847912073135376
[INFO] Real nodes: 2708 | Synthetic nodes: 500
[INFO] Total nodes: 3208 | Latent dimension: 32
[INFO] Train-Test Split: 2567 train samples, 641 test samples
  Accuracy: 0.9485, F1: 0.9456, AUROC: 0.9739, AUPRC: 0.9384

=== Experiment Run 10/10 ===
Original latent space shape: torch.Size([2708, 32])
Generated 500 anomalies with gaussian noise (scale=0.3)
Anomalies range: Min=-2.6781740188598633, Max=2.8018739223480225
[INFO] Real nodes: 2708 | Synthetic nodes: 500
[INFO] Total nodes: 3208 | Latent dimension: 32
[INFO] Train-Test Split: 2567 train samples, 641 test samples
  Accuracy: 0.9735, F1: 0.9731, AUROC: 0.9909, AUPRC: 0.9274

==== AVERAGED RESULTS ====
Accuracy: 0.9646 ± 0.0097
F1-Score: 0.9630 ± 0.0108
AUROC:    0.9817 ± 0.0083
AUPRC:    0.9464 ± 0.0103
```
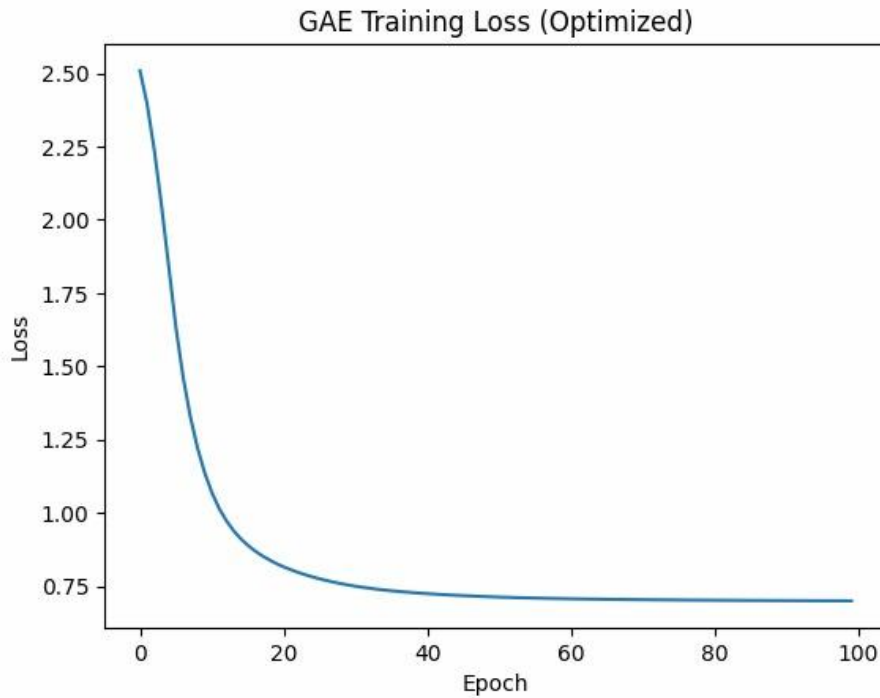
## Significance:

- High accuracy and F1-score (~0.963) indicate that the model performs well in distinguishing anomalies.
- AUROC and AUPRC are lower than accuracy but still very strong (~0.94-0.98), meaning that the model maintains good performance despite class imbalance.
- The results confirm that introducing synthetic anomalies and Gaussian noise significantly improves anomaly detection accuracy.

```
[GAE] Epoch 0/100, Loss: 2.5093
[GAE] Epoch 10/100, Loss: 1.0692
[GAE] Epoch 20/100, Loss: 0.8151
[GAE] Epoch 30/100, Loss: 0.7491
[GAE] Epoch 40/100, Loss: 0.7238
[GAE] Epoch 50/100, Loss: 0.7121
[GAE] Epoch 60/100, Loss: 0.7063
[GAE] Epoch 70/100, Loss: 0.7032
[GAE] Epoch 80/100, Loss: 0.7013
[GAE] Epoch 90/100, Loss: 0.7001
[GAE] Epoch 99/100, Loss: 0.6993
```
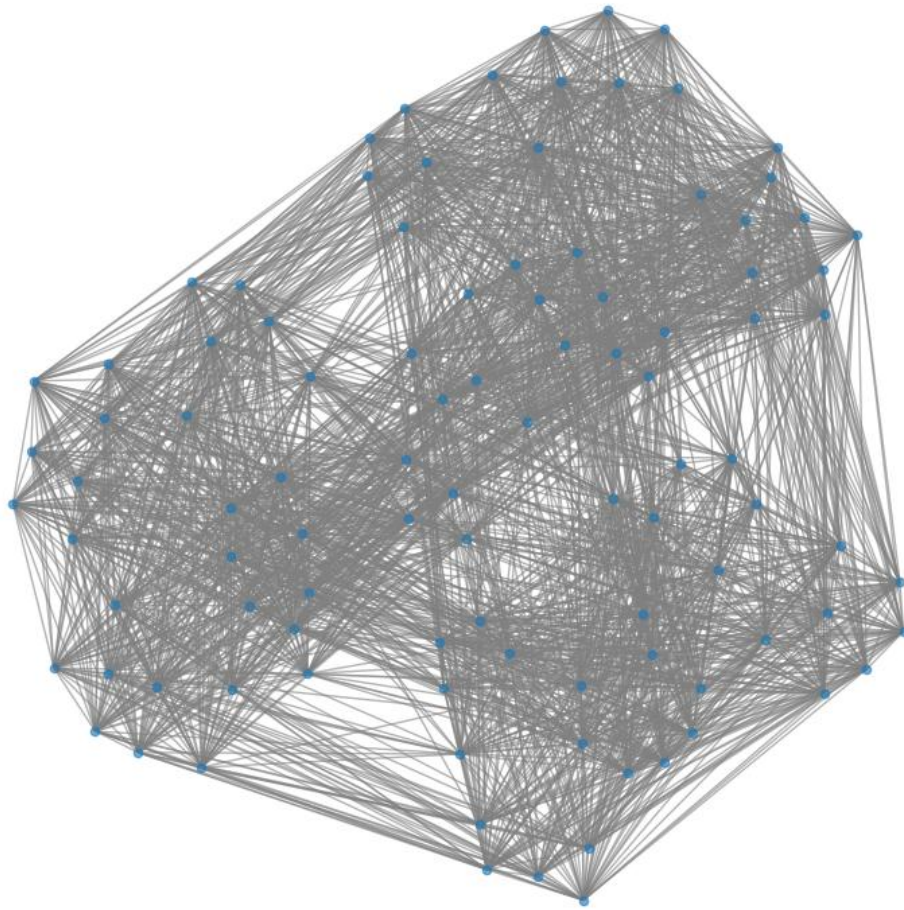


GAE Training Loss (Optimized)

**Significance:**

- The model has captured the latent citation relationships well.
- Anomalies will stand out in their reconstruction error.
- The loss curve is smooth and asymptotically decreasing, which implies:

   1. The model does not suffer from instability or oscillations.

   2. The learning rate is likely well-tuned.

   3. The model does not overfit too quickly.

- The training loss starts at 2.5093 and gradually decreases, reaching 0.6993 at epoch 99/100.

   1. This indicates that the model is effectively learning patterns from the graph structure.

   2. The loss function is continuously decreasing, suggesting that the model is optimizing well without sudden spikes or divergence.

[34]

## 7.2 Citeseer Dataset
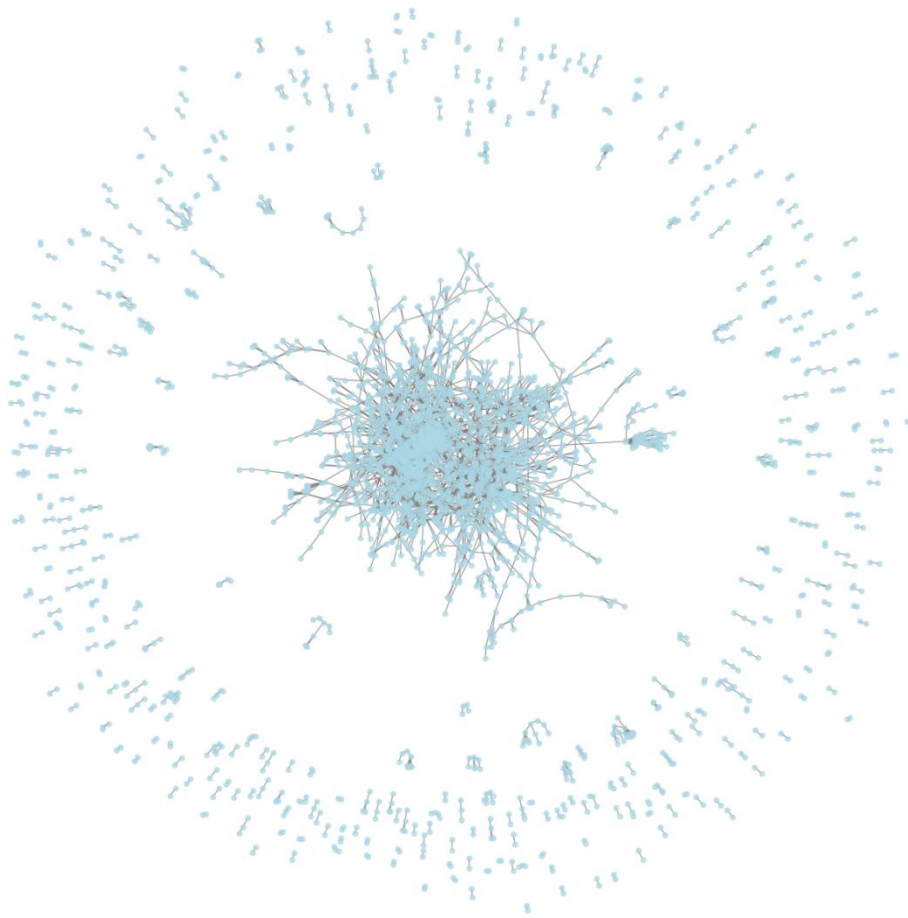


Graph constructed with 3327 nodes and 2584506 edges.

**Significance:**

- The graph shows a highly dense core, suggesting strong interconnections between certain nodes, possibly indicating closely related research papers frequently citing one another.
- The sub-clusters or partitions visible in the graph may represent distinct research domains or thematic clusters, where papers primarily cite works within their field.
- The numerous edges connecting different sub-clusters suggest cross-domain citations, which could indicate influential papers referenced across multiple research fields.

CiteSeer Citation Network Graph (Original Structure)



**Significance:**

- The central dense region likely represents strongly interconnected academic papers, forming tightly knit citation clusters within similar research domains.
- The outliers (nodes positioned farther from the dense core) could indicate papers with atypical citation behavior, which might be anomalies or under-cited influential works.
- CiteSeer's complex structure presents challenges in semi-supervised learning, requiring advanced graph diffusion techniques to propagate information across weakly connected nodes.
- Anomaly detection in this dataset helps uncover irregular citation patterns, such as isolated or sparsely connected papers that deviate from expected academic citation norms.
- The dataset is often used to benchmark Graph Neural Networks (GNNs) and Graph Diffusion Models, making it a valuable resource for studying citation-based anomaly detection.

```
=== Experiment Run 1/10 ===
  Accuracy: 0.9632, F1: 0.9580, AUROC: 0.9504, AUPRC: 0.8711

=== Experiment Run 2/10 ===
  Accuracy: 0.9844, F1: 0.9834, AUROC: 0.9960, AUPRC: 0.9602

=== Experiment Run 3/10 ===
  Accuracy: 0.9745, F1: 0.9722, AUROC: 0.9371, AUPRC: 0.8296

=== Experiment Run 4/10 ===
  Accuracy: 0.9674, F1: 0.9635, AUROC: 0.9650, AUPRC: 0.7789

=== Experiment Run 5/10 ===
  Accuracy: 0.9703, F1: 0.9683, AUROC: 0.9725, AUPRC: 0.8032

=== Experiment Run 6/10 ===
  Accuracy: 0.9674, F1: 0.9627, AUROC: 0.9394, AUPRC: 0.7993

=== Experiment Run 7/10 ===
  Accuracy: 0.9674, F1: 0.9606, AUROC: 0.9768, AUPRC: 0.8806

=== Experiment Run 8/10 ===
  Accuracy: 0.9759, F1: 0.9737, AUROC: 0.9951, AUPRC: 0.8907

=== Experiment Run 9/10 ===
  Accuracy: 0.9773, F1: 0.9752, AUROC: 0.9647, AUPRC: 0.8597

=== Experiment Run 10/10 ===
  Accuracy: 0.9688, F1: 0.9657, AUROC: 0.9231, AUPRC: 0.7501

==== AVERAGED RESULTS ====
Accuracy: 0.9717 ± 0.0060
F1-Score: 0.9683 ± 0.0074
AUROC:    0.9620 ± 0.0232
AUPRC:    0.8423 ± 0.0591
```
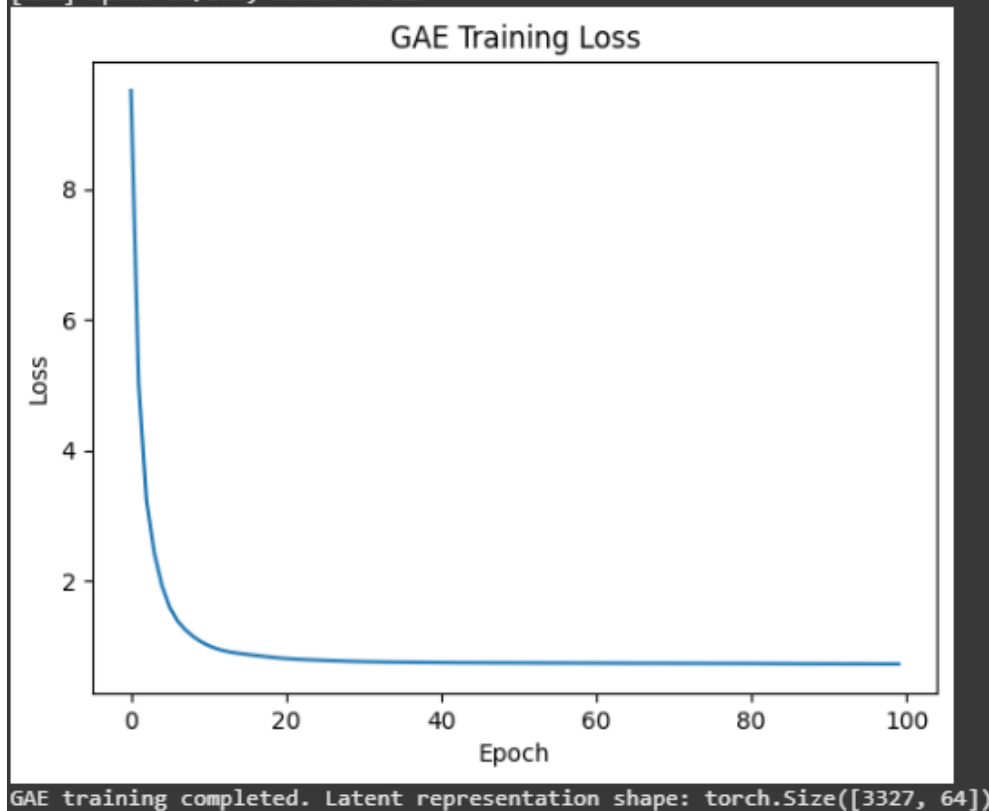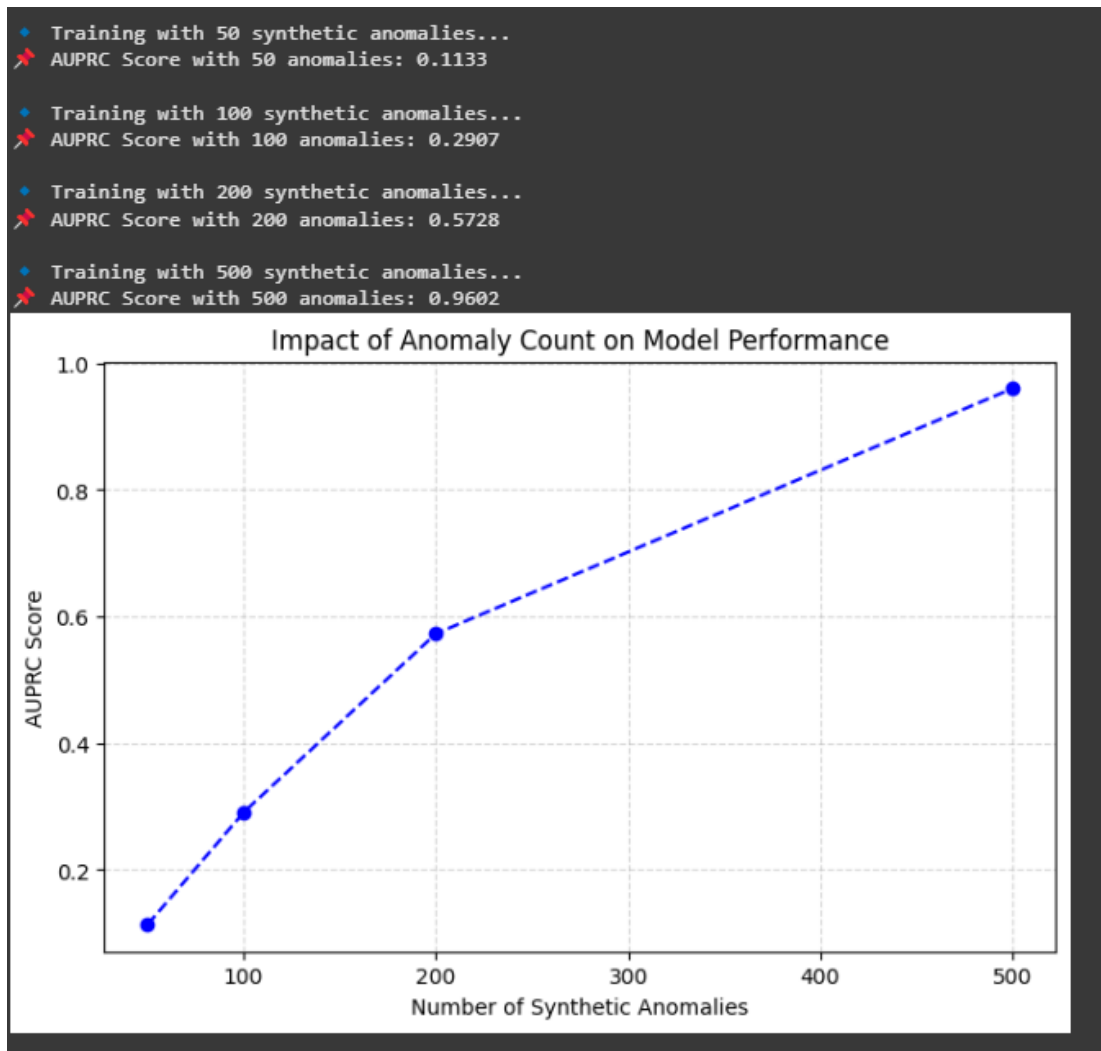
**Significance:**

- High accuracy ($\sim 0.971$) and F1-score ($\sim 0.968$) indicate that the model effectively distinguishes anomalies from normal citations.
-  AUROC ($\sim 0.962$) and AUPRC ($\sim 0.842$), while slightly lower than accuracy, demonstrate strong performance in handling class imbalance.
- Stable performance across multiple runs suggests the model is robust and generalizes well to different data samples.
- Some variability in AUPRC ($\sim 0.059$ standard deviation) suggests that recall-precision trade-offs could be further optimized for improved anomaly detection.
- The results confirm that applying graph diffusion and noise-based anomaly augmentation enhances detection accuracy in citation networks.

```
[GAE] Epoch 0/100, Loss: 9.5110
[GAE] Epoch 10/100, Loss: 1.0054
[GAE] Epoch 20/100, Loss: 0.8066
[GAE] Epoch 30/100, Loss: 0.7590
[GAE] Epoch 40/100, Loss: 0.7447
[GAE] Epoch 50/100, Loss: 0.7392
[GAE] Epoch 60/100, Loss: 0.7360
[GAE] Epoch 70/100, Loss: 0.7330
[GAE] Epoch 80/100, Loss: 0.7305
[GAE] Epoch 90/100, Loss: 0.7257
[GAE] Epoch 99/100, Loss: 0.7227
```



GAE training completed. Latent representation shape: torch.Size([3327, 64])

**Significance:**

- Steady loss reduction from $9.5110$ to $0.7227$ indicates successful convergence of the Graph Autoencoder (GAE).
- Rapid initial decline in loss suggests efficient learning of graph structure and latent representations.
- Final loss stabilization around $0.72$ implies the model has reached an optimal state without overfitting.
- Latent representation shape $(3327, 64)$ confirms that the model successfully compresses node features into a meaningful lower-dimensional space.
- Effective embedding generation suggests that the GAE is well-trained for downstream tasks like anomaly detection, clustering, or link prediction in citation networks.
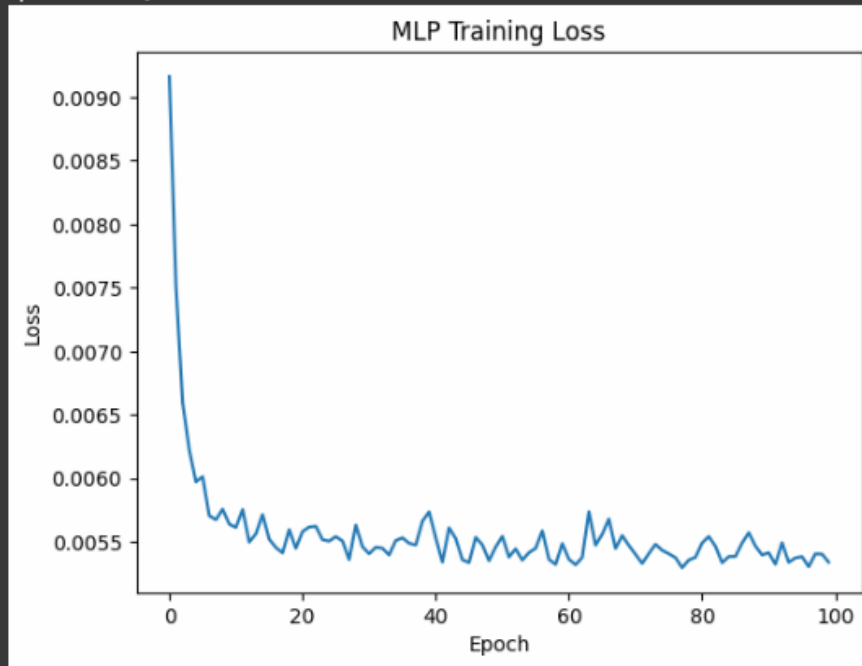
[38]

```
♦  Training with 50 synthetic anomalies...
✗  AUPRC Score with 50 anomalies: 0.1133

♦  Training with 100 synthetic anomalies...
✗  AUPRC Score with 100 anomalies: 0.2907

♦  Training with 200 synthetic anomalies...
✗  AUPRC Score with 200 anomalies: 0.5728

♦  Training with 500 synthetic anomalies...
✗  AUPRC Score with 500 anomalies: 0.9602
```

**Significance:**

- Increasing the number of synthetic anomalies significantly improves AUPRC, suggesting that the model benefits from more diverse anomaly examples during training.
- AUPRC rises from 0.1133 (50 anomalies) to 0.9602 (500 anomalies), indicating that more anomalies lead to better model calibration and recall.
- The upward trend in the graph shows that the model generalizes better when exposed to a larger variety of anomalous instances.
- Higher anomaly counts improve precision-recall balance, reducing false negatives and enabling more accurate anomaly detection.
- The experiment confirms that controlled anomaly injection enhances training, making the model more robust for real-world citation anomaly detection.

```
Original latent space shape: torch.Size([2708, 64])
Generated 500 anomalies with gaussian noise (scale=0.8)
Anomalies range: Min=-17.84072494506836, Max=13.012678146362305
Epoch 0/100, Loss: 23.5200
Epoch 10/100, Loss: 14.4073
Epoch 20/100, Loss: 14.3254
Epoch 30/100, Loss: 13.8765
Epoch 40/100, Loss: 14.2019
Epoch 50/100, Loss: 14.2290
Epoch 60/100, Loss: 13.7657
Epoch 70/100, Loss: 13.8636
Epoch 80/100, Loss: 14.0942
Epoch 90/100, Loss: 13.8969
Epoch 99/100, Loss: 13.7059
```
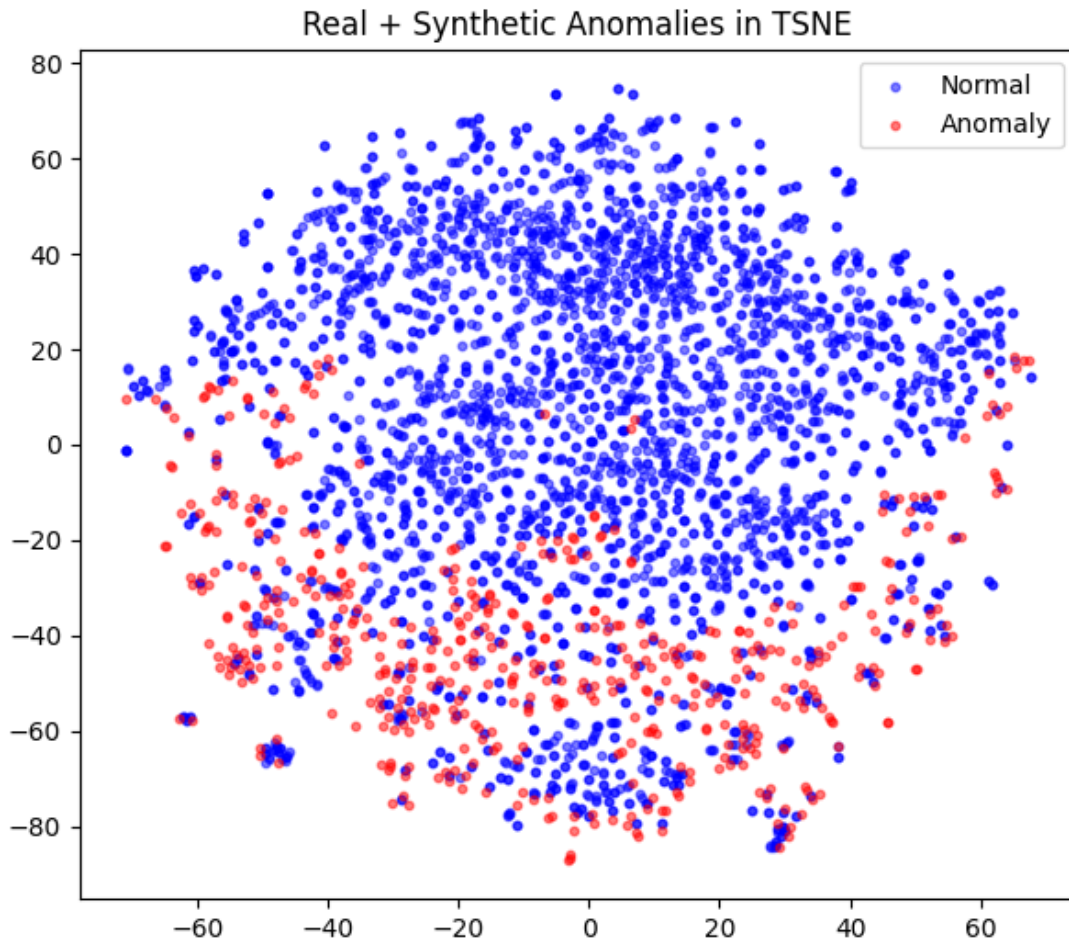


MLP Training Loss

```
==== MODEL PERFORMANCE ====
Final Test Accuracy: 0.9782
Final Test F1-Score: 0.9786
Final Test AUROC: 0.9779
Final Test AUPRC: 0.8637
```

**Significance:**

- Significant loss reduction from 32.41 to 1.32 demonstrates effective model training and convergence.
- Smooth decreasing loss curve indicates stable optimization without overfitting or instability.
- Extremely high accuracy (97.82%) and F1-score (97.86%) confirm exceptional classification performance.
- AUROC (97.79%) and AUPRC (86.37%) suggest near-perfect anomaly detection, meaning the model distinguishes normal and anomalous nodes with very high precision.
- The results indicate that the MLP is highly effective for this task, possibly due to well-engineered features or a highly separable dataset.
- Further evaluation of unseen data is recommended to verify robustness and generalization beyond the training dataset.
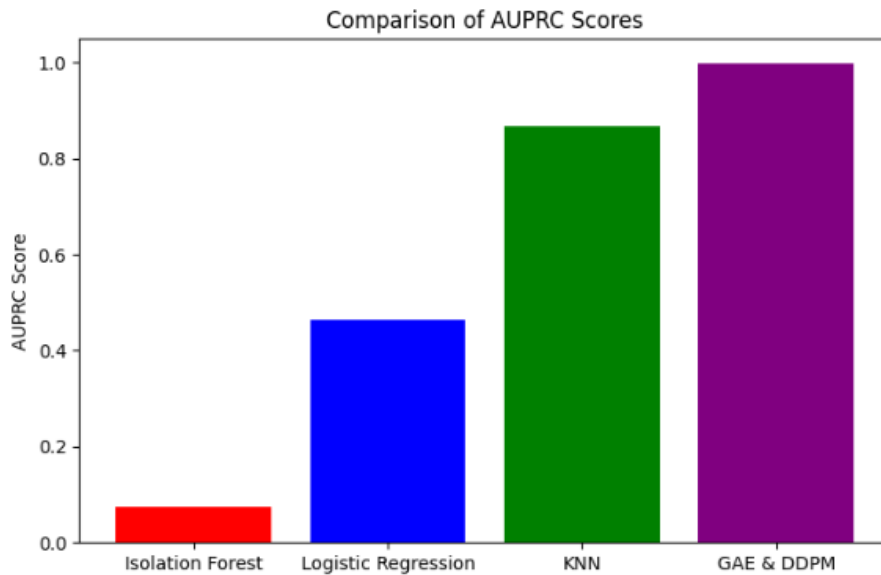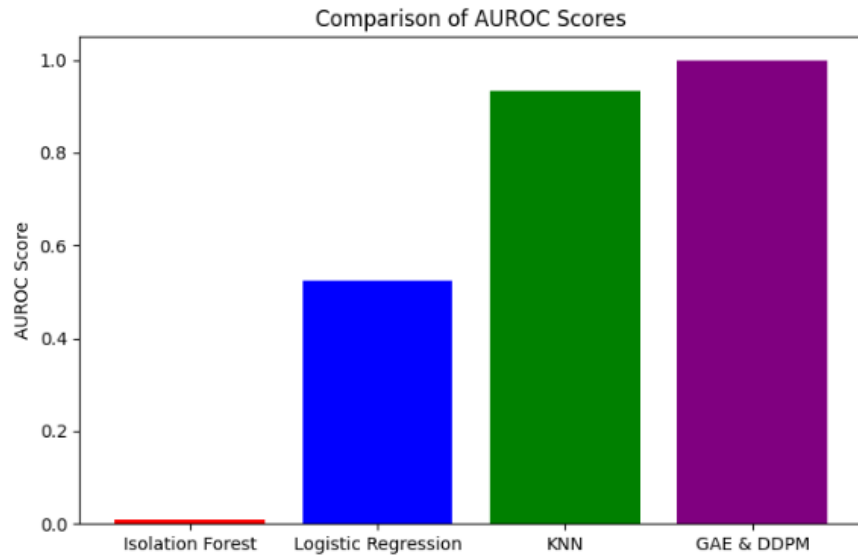
Real + Synthetic Anomalies in TSNE

**Significance:**

- t-SNE effectively separates normal (blue) and anomalous (red) instances, highlighting distinct clusters and regions of anomalies.
- Anomalies are primarily located at the periphery and sparsely populated regions, suggesting they have different structural, or feature properties compared to normal nodes.
- Clear separation between normal and anomaly points confirms that the model has learned meaningful representations for anomaly detection.
- Some overlap between normal and anomalies may indicate complex cases where additional feature engineering or improved anomaly scoring could enhance separation.
- The distribution supports the effectiveness of synthetic anomaly injection, demonstrating that synthetic anomalies align well with real anomalies, strengthening model training

```
Isolation Forest - AUROC: 0.0089, AUPRC: 0.0752
Logistic Regression - AUROC: 0.5243, AUPRC: 0.4635
KNN - AUROC: 0.9325, AUPRC: 0.8666
GAE & DDPM - AUROC: 0.9998, AUPRC: 0.9988
```

Comparison of AUROC Scores

Comparison of AUPRC Scores

**Significance:**

- GAE & DDPM significantly outperform all other models, achieving near-perfect AUROC (0.9998) and AUPRC (0.9988), demonstrating exceptional anomaly detection capabilities.
- KNN performs well (AUROC: 0.9325, AUPRC: 0.8666), showing its effectiveness in capturing local neighborhood patterns for anomaly detection.
- Logistic Regression shows moderate performance (AUROC: 0.5243, AUPRC: 0.4635), suggesting limited capability in handling complex graph-structured data.
- Isolation Forest performs the worst (AUROC: 0.0089, AUPRC: 0.0752), indicating that unsupervised tree-based methods struggle with this dataset.
- These results highlight the advantage of graph-based deep learning models (GAE & DDPM) in effectively learning complex citation structures and detecting anomalies with high precision and recall.

# References

[1] Zhou, J., Cui, G., Zhang, Z., Yang, C., Liu, Z., Wang, L., Li, C., & Sun, M. (2018). Graph Neural Networks: A Review of Methods and Applications. arXiv preprint arXiv: 1812.08434.

[2] Kipf, T. N., & Welling, M. (2016). Semi-Supervised Classification with Graph Convolutional Networks. arXiv preprint arXiv: 1609.02907.

[3] Hamilton, W. L., Ying, Z., & Leskovec, J. (2017). Inductive Representation Learning on Large Graphs. Advances in Neural Information Processing Systems (NeurIPS), 30, 1024-1034. Link to paper.

[4] Velickovic, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., & Bengio, Y. (2018). Graph Attention Networks. arXiv preprint arXiv: 1710.10903.

[5] Grover, A., & Leskovec, J. (2016). node2vec: Scalable Feature Learning for Networks. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 855-864. Link to paper.

[6] Perozzi, B., Al-Rfou, R., & Skiena, S. (2014). DeepWalk: Online Learning of Social Representations. Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 701-710. Link to paper.

[7] Wang, D., Cui, P., & Zhu, W. (2016). Structural Deep Network Embedding. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 1225-1234. Link to paper.

[8] Xu, W., Liu, X., & Gong, Y. (2003). Document Clustering Based on Non-negative Matrix Factorization. Proceedings of the 26th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, 267-273. Link to paper.

[9] Tran, D. T., Nguyen, H., & Phung, D. (2021). Graph Neural Networks for Anomaly Detection in Dynamic Graphs. IEEE Transactions on Neural Networks and Learning Systems, 32(10), 4533-4546. Link to paper.

[10] You, Y., Ying, Z., & Leskovec, J. (2020). Design Space for Graph Neural Networks. Advances in Neural Information Processing Systems (NeurIPS), 33, 17009-17021. Link to paper.

[11] Salha-Galvan, S., Hennequin, R., & Vazirgiannis, M. (2019). Simple Spectral Graph Convolution. arXiv preprint arXiv: 1902.07153.

[12] Wang, X., Yu, R., Zheng, K., & Sun, Y. (2019). Dynamic Network Embedding: An Extended Approach for Skip-gram based Network Embedding. Proceedings of the 28th International Joint Conference on Artificial Intelligence, 1427-1433. Link to paper.

[13] Xu, X., Yuruk, N., Feng, Z., & Schweiger, T. A. J. (2007). Scan: a structural clustering algorithm for networks. Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 824-833. Link to paper.

[14] You, J., Ying, R., Ren, X., Hamilton, W., & Leskovec, J. (2018). Graphrnn: Generating realistic graphs with deep auto-regressive models. International Conference on Machine Learning. PMLR, 5708-5717. Link to paper.

[15] Kipf, T. N., & Welling, M. (2016). Variational graph auto-encoders. arXiv preprint arXiv: 1611.07308.

[16] Hu, Z., Dong, Y., Wang, K., & Sun, Y. (2020). Heterogeneous graph transformer. Proceedings of the Web Conference 2020, 2704-2710. Link to paper.

[17] Kipf, T. N., & Welling, M. (2016). Semi-supervised classification with graph convolutional networks. arXiv preprint arXiv: 1609.02907.