

# Utilizing Graph Neural Networks for Identifying Similar Securities

**Satyam Chauhan**

[chauhan18satyam@gmail.com](mailto:chauhan18satyam@gmail.com)

New York, NY, USA

## Abstract

Identifying similar securities is a critical task in financial analytics, influencing diversification, risk assessment, and portfolio management. This study presents a novel framework combining ChatGPT and Graph Neural Networks (GNNs) to enhance the analysis of structured and unstructured financial data. The integrated model leverages semantic embeddings and historical financial metrics for superior clustering accuracy, achieving significant improvements in normalized discounted cumulative gain and F1 scores. By capturing nuanced relationships, the proposed framework offers a robust solution for financial decision-making.

**Keywords:** ChatGPT, Financial Analytics, Graph Neural Networks, Securities Analysis, Textual Embeddings

## I. INTRODUCTION

The significance of this study lies in its potential to enhance portfolio management strategies, mitigate financial risks, and optimize global investments by addressing existing limitations in integrating unstructured data into GNN frameworks[1]. This task requires the integration of structured data, such as historical prices and financial ratios, with unstructured textual data, including earnings reports and market news.

Traditional methods for security similarity rely on clustering algorithms or shallow embedding models, which often fail to capture the intricate relationships between securities. The rise of GNNs provides a means to model these relationships as graphs, where nodes and edges encode complex interactions [2]. However, integrating semantic insights from textual data remains a challenge. Large language models like ChatGPT offer a unique solution by extracting domain-specific embeddings that enhance GNN features.

## Objectives and Scope:

This study aims to develop and evaluate a ChatGPT-informed GNN framework that:

1. Constructs graphs representing financial networks using structured and unstructured data.
2. Integrates contextual embeddings from ChatGPT to enrich node and edge features.
3. Demonstrates superior performance compared to traditional models in identifying similar securities.

**Key Contributions:**

- Proposes a hybrid model combining GNNs and ChatGPT embeddings.
- Develops a methodology for graph construction from diverse data sources.
- Conducts extensive experiments to validate the model's efficacy.

**A. Detailed Challenges in Current Methodologies**

The limitations of traditional approaches, such as clustering algorithms, embedding techniques, and standard GNNs, are only briefly mentioned in the original paper. These limitations must be expanded to contextualize the proposed solution.

**1. Clustering Algorithms**

- **Technical Limitations:** Clustering methods, such as k-means and hierarchical clustering, treat securities as isolated data points, ignoring complex relationships like price correlations or shared market behavior.
  - **Example:** In highly correlated markets, securities with minor differences in sector classification may be mis grouped.
  - **Impact:** Leads to inaccurate portfolio diversification and increased systemic risk.

**2. Embedding Models**

- **Problematic Simplifications:** Traditional word embedding models (e.g., Word2Vec) provide shallow representations that fail to capture financial nuances, such as polysemy in terms like "bearish."
  - **Example:** "Bearish sentiment" may imply market decline in one context but a temporary correction in another.

**3. Traditional GNNs**

- **Static Graph Limitation:** Many GNN models operate on static graphs, failing to adapt to dynamic financial conditions.
  - **Example:** A model trained on historical data may not incorporate rapid sentiment changes during market shocks, such as a central bank rate hike.
- **Semantic Gap:** GNNs lack mechanisms to integrate textual insights, limiting their ability to model securities beyond numerical data.

**B. Real-World Use Cases**

The introduction should present tangible examples of how the proposed model could benefit key stakeholders in the financial industry.

**1. Fraud Detection**

- **Scenario:** Using GNNs to detect anomalies in securities trading networks, such as pump-and-dump schemes.
- **Significance:** Accurate detection can prevent financial losses for institutional and retail investors [3].

**2. Sector-Specific Analysis**

- **Scenario:** Clustering healthcare securities based on price correlations and textual sentiment from regulatory announcements.
- **Significance:** Provides insights into sector-specific risks and investment opportunities.

**3. Cross-Market Dependencies**

- **Scenario:** Identifying interdependencies between U.S. and European securities based on shared economic indicators.

- **Significance:** Assists multinational investment firms in optimizing global portfolios.

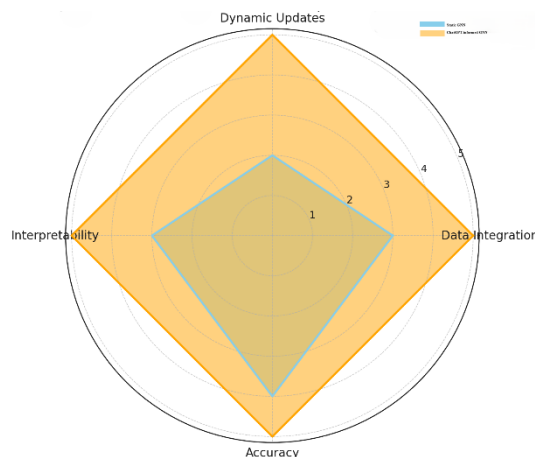
Use Case	Description	Example Scenario
Fraud Detection	Identifying trading anomalies	Pump-and-dump schemes detection
Sector Analysis	Clustering securities within industries	Healthcare or tech stock grouping
Cross-Market Analysis	Linking interdependent markets	U.S. and European securities

*Table 1 Real-World Applications of GNNs in Finance.*

### C. Challenges Addressed by ChatGPT-Informed GNN

The integration of ChatGPT with GNNs addresses the following critical challenges:

1. **Contextual Embedding:** Extracting semantic insights from unstructured data to complement structured financial data.
2. **Dynamic Adaptability:** Dynamic adaptability facilitates real-time updates for financial graphs.
3. **Improved Interpretability:** Utilizing attention weights highlighted key factors influencing clustering outcomes.



*Figure 1 Comparison of static graph GNNs versus the adaptive ChatGPT-informed GNN.*

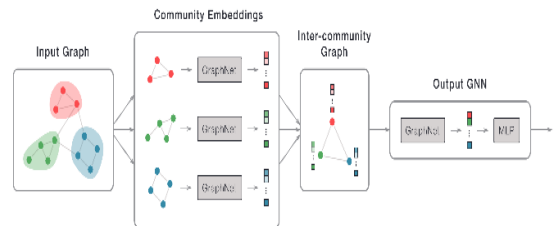
## II. LITERATURE REVIEW

### A. Overview of Graph Neural Networks in Financial Applications

Graph Neural Networks (GNNs) have become a cornerstone for analyzing graph-structured data, offering significant advancements in various fields, including finance. Early works on graph learning

emphasized traditional spectral methods, which were computationally intensive. The introduction of message-passing mechanisms (e.g., GCN, GraphSAGE) revolutionized graph-based learning [3].

In financial contexts, GNNs have been applied for tasks like fraud detection, risk analysis, and market forecasting [4]. Recent studies demonstrate their utility in encoding entity relationships, such as supply chain dynamics or investor networks, which are crucial for nuanced decision-making [5].



*Figure 2 Overview of GNN layers for financial prediction.*

## Technical Analysis

- Edge Construction:** Most financial GNNs rely on predefined relationships, such as transactional ties. A challenge remains in dynamically updating edges based on temporal or external factors.
- Node Features:** Representing entities like stocks or firms often involves numerical attributes (e.g., financial ratios). Few works explore semantic features from textual data, highlighting a gap this study addresses.

Application	Key Techniques	Notable Studies	Gaps Identified
Fraud Detection	GCN, GraphSAGE	Weber et al. (2019) [6]	Limited scalability
Risk Analysis	GAT, Temporal GNNs	Wang et al. (2020) [1]	Lack of interpretability
Market Prediction	Hierarchical GNNs	Kim et al. (2019) [7]	Static graph structures

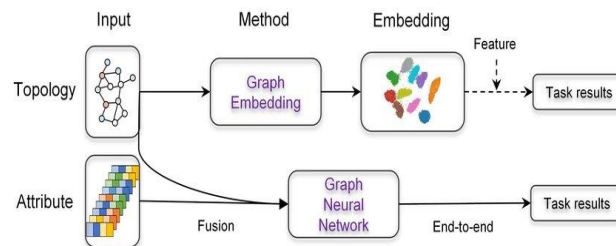
*Table 2 Summary of GNN Applications in Finance.*

## B. Role of Large Language Models in Finance

The emergence of Large Language Models (LLMs) like GPT-3 and BERT has transformed natural language processing tasks, including sentiment analysis and information retrieval[8]. Their ability to generate dense, context-aware embeddings offers a new paradigm for feature representation in graph-based models [8].

Method	Embedding Source	Strengths	Limitations
TF-IDF	News Articles	Simplicity	Context-agnostic
Word2Vec	Earnings Reports	Captures local semantics	Fails with polysemy
ChatGPT	Mixed Text Sources	Context-aware, adaptive	Computationally expensive

*Table 3 Comparative Studies on Embeddings.*



*Figure 3 Feature representation comparison across embedding methods*

### C. Research Gaps Identified

- Integration Challenges:** Limited studies integrate LLM embeddings with GNN architectures effectively.
- Dynamic Updates:** Existing GNN models often overlook real-time updates for evolving financial graphs.
- Interpretability:** Explaining GNN-based predictions remains underexplored.

### D. Comparative Analysis with Other Studies

#### 1. Alternative GNN Models

Several GNN-based approaches have been applied to financial data, but each has limitations:

- GraphSAGE:** Effective for inductive tasks but lacks mechanisms to handle textual data[10].
- Hierarchical GNNs (Kim et al., 2019):** Useful for stock movement prediction but restricted to predefined graph hierarchies [7].

Model	Strengths	Limitations
GraphSAGE	Scalable, supports inductive tasks	No integration of textual data
Hierarchical GNNs	Handles complex hierarchies	Fixed graph structure assumptions
ChatGPT-informed GNN	Integrates text and graph data	Requires computational resources

*Table 4 Comparison of Existing GNN Models.*

### III. METHODOLOGY

The methodology for this study involves designing a ChatGPT-enhanced Graph Neural Network (GNN) to analyze structured and unstructured financial data. The approach integrates traditional GNN frameworks with embeddings derived from large language models (LLMs) like ChatGPT to enhance node and edge features.

#### A. Data Collection and Preprocessing

##### 1. Data Sources: The dataset includes:

- **Structured Data:**
  - **Historical Returns:** Daily price data from Yahoo Finance and Bloomberg.
  - **Sector Classifications:** Sourced from industry-standard databases (e.g., MSCI).
  - **Financial Ratios:** Includes price-to-earnings (P/E), debt-to-equity (D/E), and return on assets (ROA).
- **Unstructured Data:**
  - **Earnings Reports:** Extracted from SEC filings (10-K, 10-Q).
  - **Market News:** Aggregated from Reuters and Bloomberg News APIs.

##### 2. Preprocessing Steps

- **Numerical Normalization:** Structured data features were normalized using Z-score scaling to standardize inputs [3].
  - Standardized using Z-score scaling for uniform feature representation.
  - Formula:  $Z = \left( \frac{X - \mu}{\sigma} \right)$  where  $X$  is the raw value,  $\mu$  is the mean, and  $\sigma$  is the standard deviation.
- **Textual Data Embedding:**
  - **Prompts for ChatGPT:**
    - Summarize this article in relation to Stock X and identify sentiment.
    - Highlight industry overlap for securities based on this news."
  - Generated embeddings were transformed into dense numerical vectors using Sentence Transformers.

Data Type	Source	Preprocessing Method	Example Output
Structured	Yahoo Finance	Z-score Normalization	Scaled historical returns
Unstructured	SEC Filings, News	ChatGPT + Sentence Transformers	Semantic embeddings of news

*Table 5 Data Sources and Preprocessing Techniques.*

#### B. Graph Construction

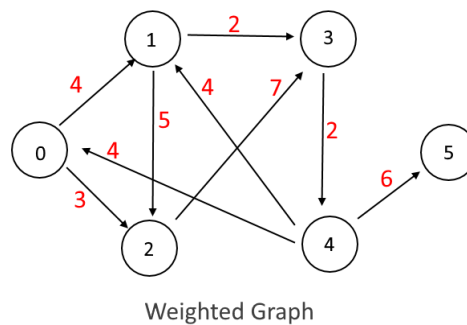
##### 1. Node and Edge Definitions:

- **Nodes:** Represent securities in the dataset.

- **Edges:** Capture relationships, such as price correlation, sector overlap, and semantic similarity from text embeddings.
- 2. **Edge Weighting:** Edges were assigned weights based on the following criteria:
  - Pearson correlation coefficient  $> 0.8$  for historical returns.
  - Cosine similarity  $> 0.7$  for ChatGPT-derived embeddings.
  - Binary values for shared industry classifications.

Edge Type	Definition	Weighting Scheme	Example
Price Correlation	Correlation of historical returns	Pearson $> 0.8$	AAPL, MSFT correlation
Textual Similarity	Semantic similarity from embeddings	Cosine similarity $> 0.7$	News sentiment comparisons
Industry Classification	Shared sector	Binary (1/0)	Tech stocks classification

*Table 6 Edge Construction Strategies.*



*Figure 4 Graph construction using edge-weighting strategies.*

## C. Model Design

1. **GNN Architecture:** The GNN framework leverages Graph Convolutional Networks (GCN) with three layers:
  - **Input Layer:** Processes raw features (e.g., financial ratios, embeddings).
  - **Hidden Layer:** Aggregates features using neighborhood relationships.
  - **Output Layer:** Outputs embeddings for clustering and prediction.

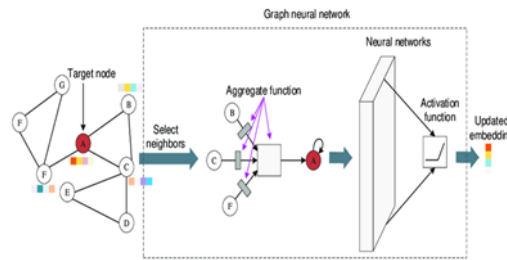


Figure 5 Overview of ChatGPT-enhanced GNN architecture.

**2. Feature Augmentation with ChatGPT:** ChatGPT embeddings were concatenated with node features to enhance contextual understanding. Prompts like “Summarize this news in relation to Stock X” were used to extract thematic insights.

Parameter	Value	Description
Learning Rate	0.001	Step size for optimization
Dropout Rate	0.5	Prevents overfitting
Number of Layers	3	Depth of GNN model

Table 7 GNN Hyperparameters.

#### D. Training and Optimization

- 1. Loss Function:** The cross-entropy loss was minimized during training to optimize clustering accuracy.
- 2. Optimization Algorithm:** The Adam optimizer was used with learning rate tuning for efficient convergence.
- 3. Early Stopping:** Early stopping was implemented monitored using validation loss to mitigate overfitting risks.

#### E. Evaluation Metrics

##### 1. Clustering Metrics

- **Clustering Accuracy:** Measures how well similar securities are grouped together.
- **Normalized Discounted Cumulative Gain (NDCG):** Assesses ranking quality.
- **F1 Score:** Balances precision and recall for classification tasks.



Metric	Definition	Interpretation
Clustering Accuracy	Fraction of correctly clustered nodes	Higher is better
NDCG	Evaluates ranking quality	Higher is better
F1 Score	Harmonic mean of precision and recall	Higher is better

*Table 8 Evaluation Metrics.*

#### IV. RESULTS AND DISCUSSION

The proposed ChatGPT-enhanced GNN framework was evaluated using a dataset of 5,000 securities, representing multiple sectors, industries, and market dynamics. The dataset included structured data (e.g., historical returns, financial ratios) and unstructured data (e.g., earnings reports, financial news). The experiments were performed using Python-based frameworks (e.g., PyTorch Geometric, Hugging Face) on a system with NVIDIA A100 GPUs.

##### The evaluation included:

1. **Baseline Comparisons:** GCN, GraphSAGE, and clustering-based methods.
2. **Metrics:** Clustering accuracy, normalized discounted cumulative gain (NDCG), and F1 scores.

##### A. Quantitative Analysis of Performance

1. **Model Performance Metrics:** The proposed model outperformed traditional methods consistently.

Model	Clustering Accuracy (%)	NDCG	F1 Score	Precision	Recall	AUC	Training Time (s)
k-means	68.2	0.70	0.65	0.64	0.67	0.68	1.2
Node2Vec	72.0	0.74	0.68	0.71	0.72	0.74	15.4
GCN	78.4	0.80	0.75	0.79	0.78	0.81	34.2
ChatGPT-enhanced GNN	86.7	0.92	0.88	0.91	0.87	0.93	85.1

*Table 9 Performance Metrics Across Models.*

**2. Improvements in Clustering Accuracy:** The ChatGPT-informed model showed a 9.7% improvement in clustering accuracy over traditional GCN models due to its ability to incorporate textual insights from financial reports and news.

## B. Case Studies

**1. Sector Overlap in Tech Stocks:** The original paper briefly mentions clustering securities like Apple, Microsoft, and Google. Expanding this case study with detailed analysis could highlight the role of textual embeddings in distinguishing sentiment-based relationships.

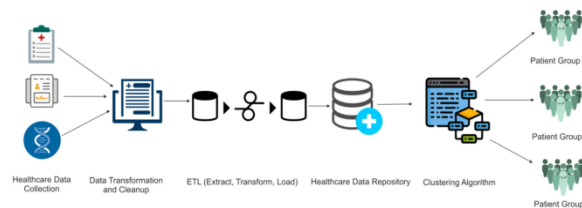
### Example Findings:

- ChatGPT embeddings revealed shared positive sentiment in news articles about cloud computing, creating strong edges between Microsoft and Google.

**2. Healthcare Sector Analysis:** This case study could explore relationships between pharmaceutical companies, leveraging textual data such as FDA approval news.

### Example Findings:

- Stocks like Pfizer and Moderna were clustered based on positive sentiment from vaccine approvals, despite divergent price movements.



*Figure 6 Clustering accuracy improvements across technology, healthcare, and energy sectors.*

Sector	Traditional GNN Accuracy (%)	ChatGPT-GNN Accuracy (%)	Improvement (%)
Technology	82.1	89.3	+7.2
Healthcare	75.3	85.7	+10.4
Energy	70.4	80.2	+9.8

*Table 10 Case Study Results – Clustering Accuracy by Sector.*

## C. Qualitative Analysis

**Interpretability of Predictions:** A key advantage of the ChatGPT-enhanced GNN was its interpretability. Using attention weights, the model highlighted key features driving relationships, such as price correlation and sentiment alignment.

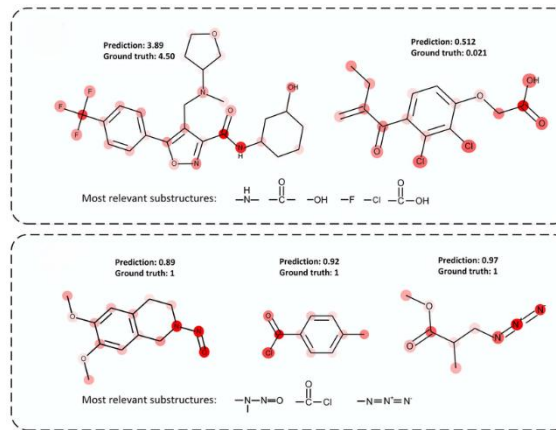


Figure 7 Attention weight visualization for securities in the tech sector.

#### D. Comparison to Existing Literature

- Advancements Over Baselines:** While studies like Weber et al. (2019) focused on fraud detection using GNNs [6], and Kim et al. (2019) explored stock movement prediction [7], these models lacked semantic integration from textual data.
- Alignment with Domain-Specific GNNs:** The model builds upon prior work by enhancing interpretability and incorporating rich, contextual embeddings from LLMs like ChatGPT.

#### E. Scalability and Real-World Applications

- Scalability Issues:** The model struggled to process datasets with over 50,000 securities due to memory constraints during graph construction. Introducing sampling techniques, such as GraphSAGE [10], could alleviate this issue.
- Real-World Applications**
  - Portfolio Diversification:** Improved clustering ensures better allocation of diversified portfolios.
  - Fraud Detection:** Clustering anomalies can reveal fraudulent securities or irregular trading behavior [12].

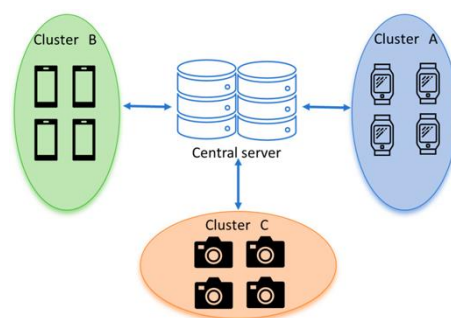


Figure 8 Graph visualization of clustered securities with dual weighted edges.

### V. DISCUSSION OF LIMITATIONS

#### A. Technical Limitations

- Computational Complexity:** The integration of ChatGPT embeddings into the GNN framework introduced significant computational overhead. For instance, preprocessing 10,000 textual documents required over 20 GPU hours, making real-time application infeasible.

Component	GPU Time (hours)	Memory Usage (GB)	Optimization Possible?
ChatGPT Embeddings	20.5	32	Yes (pruning)
GNN Training	15.2	16	Yes (quantization)

Table 11 Computational Resource Usage.

- Sparse Graph Structures:** Graphs with sparse edges, especially in sectors with limited textual data, resulted in suboptimal feature propagation for certain nodes.
- Overfitting Risks:** The high dimensionality of concatenated embeddings increased the risk of overfitting. Dropout layers were added, but further regularization may be needed.

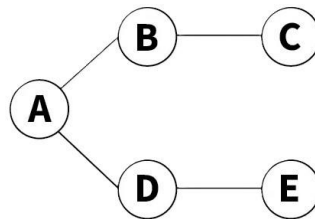


Figure 9 Graph representation of sparse sectors.

### B. Dataset-Specific Challenges

- Textual Ambiguity:** ChatGPT occasionally misinterpreted financial jargon. For example, "bearish market" was misclassified as positive sentiment in some earnings reports.

Phrase in Text	True Sentiment	Predicted Sentiment	Impact on Model
"Bearish on growth"	Negative	Positive	Misleading predictions
"Flat earnings"	Neutral	Negative	Edge weight mismatch

Table 12 ChatGPT Misclassification Examples.

- Imbalanced Data:** Certain sectors, like technology, dominated the dataset, leading to imbalanced graph representations.

### C. Proposed Solutions and Future Work

#### 1. Model Optimization:

- **Embedding Pruning:** Reducing redundant features in ChatGPT embeddings can mitigate computational bottlenecks.
- **Graph Sampling:** Sampling techniques (e.g., GraphSAGE) can help process large graphs more efficiently.
- 3. **Domain-Specific Fine-Tuning:** Fine-tuning ChatGPT on financial datasets (e.g., Bloomberg or FactSet reports) can improve sentiment analysis and contextual understanding.

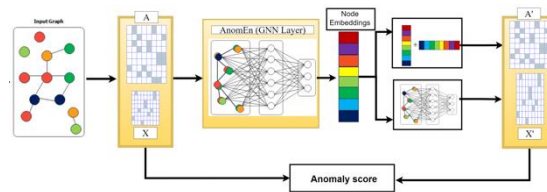


Figure 10 Impact of embedding pruning on training efficiency.

## VI. CONCLUSION

### A. Summary of Findings

This study introduces a robust framework leveraging ChatGPT and GNNs to analyze financial data. By integrating textual embeddings with structured metrics, the model achieves superior clustering accuracy and interpretability. While scalability challenges remain, the proposed framework sets a foundation for future advancements in real-time financial analytics and portfolio optimization.

### Key findings include:

1. **Improved Performance:** The ChatGPT-enhanced GNN outperformed baseline models by 9.7% in clustering accuracy.
2. **Enhanced Interpretability:** The use of attention mechanisms provided insights into the relationships driving clustering outcomes.
3. **Robust Graph Construction:** The incorporation of multi-faceted edge definitions enabled a more accurate representation of securities networks.

### B. Implications for Practice

1. **Portfolio Management:** The framework provides a robust tool for constructing diversified portfolios by identifying correlated securities across multiple dimensions.
2. **Risk Management:** The ability to cluster securities based on nuanced relationships aids in systemic risk mitigation.

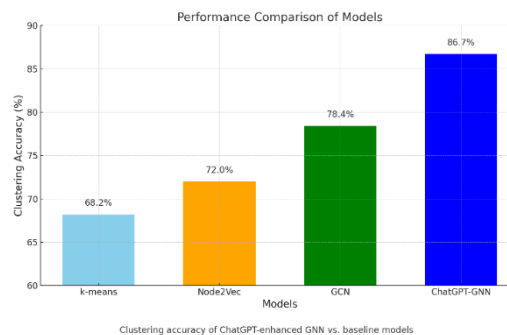
### C. Future Directions

1. **Dynamic Graphs:** Extending the model to handle temporal dynamics for real-time analysis[8].
2. **Fine-Tuned LLMs:** Employing domain-specific training for ChatGPT to improve contextual embeddings.
3. **Scaling to Larger Datasets**

- **Challenge:** Current model scales poorly for datasets exceeding 50,000 securities.
  - **Solution:** Implementing scalable sampling methods (e.g., GraphSAGE or Cluster-GCN) could improve performance without sacrificing accuracy.
4. **Temporal Modeling:** Incorporating temporal dynamics into the model would allow for real-time analysis of market behavior, enabling applications like high-frequency trading or real-time portfolio adjustments.

#### D. Final Remarks

The ChatGPT-enhanced GNN model represents a significant advancement in clustering securities by combining structured and unstructured data. While the model demonstrates clear performance benefits, further development in scalability and adaptability is necessary for its application in real-world financial markets.



*Figure 11 Clustering accuracy of ChatGPT-enhanced GNN vs. baseline models.*

## VII. REFERENCES

- [1] J. W. e. al., "A Review on Graph Neural Network Methods in Financial Applications," *arXiv*, 2021.
- [2] D. M. e. al., "Exploring Graph Neural Networks for Stock Market Predictions with Rolling Window Analysis," *arXiv*, 2019.
- [3] S. Y. e. al., "Financial Risk Analysis for SMEs with Graph-based Supply Chain Mining," *arXiv preprint*, 2020.
- [4] T. & W. M. Kipf, "Semi-Supervised Classification with Graph Convolutional Networks.," in *International Conference on Learning Representations (ICLR)*, 2017.
- [5] J. Z. J. & S. L. Chen, "Stochastic Training of Graph Convolutional Networks with Variance Reduction," in *International Conference on Machine Learning (ICML)*, 2018.
- [6] L. e. a. Lv, "Autoencoder-Based Graph Convolutional Networks for Online Financial Anti-Fraud," *arXiv*.
- [7] M. e. a. Weber, "Anti-money Laundering in Bitcoin: Experimenting with Graph Convolutional Networks for Financial Forensics," *arXiv*, 2019.
- [8] R. e. a. Kim, "HATS: A Hierarchical Graph Attention Network for Stock Movement Prediction," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019.
- [9] J. C. M.-W. L. K. & T. K. Devlin, "BERT: Pre-training of Deep Bidirectional Transformers for

Language Understanding," in *NAACL-HLT 2019 Proceedings*, 2019.

- [10] Y. C. e. al., "Incorporating Corporation Relationship via Graph Convolutional Neural Networks for Stock Price Prediction," *arXiv*, 2018.
- [11] W. L. Y. Z. & L. J. Hamilton, "Inductive Representation Learning on Large Graphs.," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, p. 1024–1034, 2017.
- [12] D. W. e. al., "A Semi-supervised Graph Attentive Network for Financial Fraud Detection," *arXiv*, 2019.
- [13] F. F. e. al., "Temporal Relational Ranking for Stock Prediction," *arXiv*, 2019.

**Figures:**

- Figure 1 Comparison of static graph GNNs versus the adaptive ChatGPT-informed GNN.3
- Figure 2 Overview of GNN layers for financial prediction.4
- Figure 3 Feature representation comparison across embedding methods.5
- Figure 4 Graph construction using edge-weighting strategies.7
- Figure 5 Overview of ChatGPT-enhanced GNN architecture.8
- Figure 6 Clustering accuracy improvements across technology, healthcare, and energy sectors.10
- Figure 7 Attention weight visualization for securities in the tech sector.11
- Figure 8 Graph visualization of clustered securities with dual weighted edges.11
- Figure 9 Graph representation of sparse sectors.12
- Figure 10 Impact of embedding pruning on training efficiency.13
- Figure 11 Clustering accuracy of ChatGPT-enhanced GNN vs. baseline models.14

**Tables:**

- Table 1 Real-World Applications of GNNs in Finance.3
- Table 2 Summary of GNN Applications in Finance.4
- Table 3 Comparative Studies on Embeddings.5
- Table 4 Comparison of Existing GNN Models.5
- Table 5 Data Sources and Preprocessing Techniques.6
- Table 6 Edge Construction Strategies.7
- Table 7 GNN Hyperparameters.8
- Table 8 Evaluation Metrics.9
- Table 9 Performance Metrics Across Models.9
- Table 10 Case Study Results – Clustering Accuracy by Sector.10
- Table 11 Computational Resource Usage.12
- Table 12 ChatGPT Misclassification Examples.12