

National Cheng Kung University

MS Degree Program on AI Robotics

Master's Thesis

生成式多視角互動圖神經網路之少樣本假新聞檢測

GemGNN: Generative Multi-view Interaction Graph Neural Networks for Few-shot
Fake News Detection

學生：余振揚

Student：Chen-Yang Yu

指導老師：李政德 博士

Advisor：Dr. Cheng-Te Li

July 2025

Abstract

Fake news has become a critical threat to information integrity and social stability, particularly in few-shot scenarios where limited labeled data is available for emerging topics or misinformation campaigns. Traditional fake news detection methods rely heavily on user propagation patterns or require extensive labeled datasets, making them impractical for real-world deployment where such data is scarce or unavailable due to privacy constraints. This thesis presents GemGNN (Generative Multi-view Interaction Graph Neural Networks), a novel framework for few-shot fake news detection that addresses these fundamental limitations through innovative heterogeneous graph neural network modeling.

Our approach introduces four key innovations that collectively establish a new paradigm for content-based misinformation detection: First, we develop a generative user interaction simulation method using Large Language Models (LLMs) to synthesize diverse user interactions with multiple semantic tones (neutral, affirmative, skeptical), effectively creating a heterogeneous node structure that captures both content and social context without requiring real user propagation data. Second, we propose a Test-Isolated K-Nearest Neighbor (KNN) edge construction strategy that prevents information leakage between test nodes during graph construction, ensuring more realistic and robust evaluation protocols in few-shot scenarios where data leakage can lead to overoptimistic performance estimates. Third, we implement a multi-view graph construction approach that partitions news embeddings into multiple semantic perspectives, enabling the capture of diverse content aspects through parallel graph structures and enhanced data augmentation at the graph level. Fourth, we design a specialized Heterogeneous Graph Attention Network (HAN) architecture that models complex type-specific relationships between news articles and generated user interactions through hierarchical attention mechanisms and meta-path based message passing.

The GemGNN framework operates under a transductive learning paradigm where all nodes (labeled, unlabeled, and test) participate in graph message passing, but only labeled nodes contribute to loss computation, maximizing the utility of limited supervision. Our heterogeneous architecture distinguishes between news nodes (characterized by rich textual embeddings) and interaction nodes (characterized by sentiment and tone features), enabling the model to learn distinct representation spaces for different entity types while capturing their relationships through learned attention weights.

Extensive experiments across comprehensive parameter grids on the FakeNewsNet datasets (PolitiFact and GossipCop) demonstrate that GemGNN significantly outperforms state-of-the-art methods across various few-shot configurations ($K=3-16$ samples per class). Our

method achieves superior performance compared to traditional approaches (MLP, LSTM), transformer-based models (BERT, RoBERTa, DeBERTa), large language models (LLaMA, Gemma), and existing graph-based methods (Less4FD, HeteroSGT, BertGCN). Comprehensive ablation studies validate the effectiveness of each architectural component, demonstrating that the synergistic combination of generative interactions, test-isolated KNN edge construction, multi-view graph decomposition, and heterogeneous attention mechanisms provides substantial and consistent improvements in few-shot fake news detection performance.

The contributions of this work establish a new paradigm for privacy-preserving fake news detection that eliminates dependency on user behavior data while maintaining superior performance in data-scarce scenarios. The framework is particularly suitable for emerging misinformation detection tasks, privacy-sensitive applications, and scenarios where social interaction data is unavailable or unreliable, addressing critical gaps in current detection capabilities for real-world deployment.

Keyword: Fake News Detection, Few Shot Learning, Transductive Learning, Generative Interaction, Graph Neural Network



Table of Contents

Abstract	i
Table of Contents	iii
List of Tables	vi
List of Figures	vii
Chapter 1. Introduction	1
1.1. Research Background and Motivation	1
1.2. Research Contributions	2
1.3. Thesis Organization	3
Chapter 2. Problem Statement	5
2.1. Basic Concepts and Notations	5
2.2. Formal Problem Definition	6
2.3. Key Challenges and Constraints	6
Chapter 3. Related Work	9
3.1. Traditional Fake News Detection Methods	9
3.1.1. Feature Engineering Approaches	9
3.1.2. Sequential Models	10
3.2. Deep Learning Approaches	10
3.2.1. Transformer-based Models	10
3.2.2. Large Language Models for Fake News Detection	11
3.3. Graph-based Fake News Detection	11
3.3.1. Document-level Graph Classification	12
3.3.2. User Propagation-based Methods	12
3.3.3. Heterogeneous Graph Neural Networks	13
3.4. Large Language Models and Graph Enhancement	14
3.4.1. LLM-Enhanced Graph Construction	14
3.4.2. LLM Direct Detection Approaches	14
3.5. Few-Shot Learning in NLP	15
3.5.1. Few-Shot Learning Fundamentals	15
3.5.2. Meta-Learning Approaches	16
3.5.3. Contrastive Learning and Data Augmentation	17
3.6. Graph Neural Networks for Fake News Detection	17
3.6.1. Message Passing Framework	17
3.6.2. Heterogeneous Graph Neural Networks	18
3.6.3. Heterogeneous Graph Attention Networks	18
3.6.4. Graph Construction Strategies	19
3.7. Limitations of Existing Methods	20

Chapter 4. Methodology: GemGNN Framework	22
4.1. Framework Overview	22
4.2. Generative User Interaction Simulation	23
4.2.1. LLM-based Interaction Generation	23
4.2.2. Multi-tone Interaction Design	23
4.2.3. Interaction-News Edge Construction	24
4.3. Graph Construction Methodologies: KNN vs Test-Isolated KNN	24
4.3.1. Traditional KNN: Performance-Optimized Graph Construction	25
4.3.2. Test-Isolated KNN: Evaluation-Realistic Graph Construction	25
4.3.3. Performance vs. Realism Trade-off Analysis	26
4.3.4. Deployment Context Decision Framework	27
4.3.5. Technical Implementation Details	27
4.4. DeBERTa vs RoBERTa: Text Encoder Selection Rationale	28
4.4.1. Disentangled Attention and Embedding Structure	28
4.4.2. Multi-View Embedding Partitioning Advantages	28
4.4.3. Empirical Validation of Encoder Choice	29
4.4.4. Computational and Practical Considerations	29
4.5. Multi-View Graph Construction	30
4.5.1. Embedding Dimension Splitting Strategy	30
4.5.2. View-specific Edge Construction	30
4.5.3. Multi-Graph Training Strategy	31
4.6. Heterogeneous Graph Architecture	31
4.6.1. Node Types and Features	31
4.6.2. Edge Types and Relations	31
4.6.3. HAN-based Message Passing and Classification	32
4.7. Loss Function Design and Training Strategy	32
4.7.1. Enhanced Loss Functions for Few-Shot Learning	32
4.7.2. Transductive Learning Framework	33
Chapter 5. Experimental Setup	34
5.1. Datasets and Preprocessing	34
5.1.1. FakeNewsNet Datasets	34
5.1.2. Data Statistics and Characteristics	35
5.1.3. Text Embedding Generation	35
5.2. Baseline Methods	36
5.2.1. Traditional Methods	36
5.2.2. Language Models	36
5.2.3. Large Language Models	36
5.2.4. Graph-based Methods	37
5.3. Evaluation Methodology	37
5.3.1. Few-Shot Evaluation Protocol	37
5.3.2. Performance Metrics	38
5.3.3. Statistical Significance Testing	38
5.4. Implementation Details	38
5.4.1. Hyperparameter Settings	38
5.4.2. Model Architecture Configuration	39
5.4.3. Training Configuration and Hardware Setup	39

Chapter 6. Results and Analysis	41
6.1. Main Results	41
6.1.1. Performance on PolitiFact Dataset	41
6.1.2. Performance on GossipCop Dataset	42
6.1.3. Comparison with Baseline Methods	43
6.2. Ablation Studies	43
6.2.1. Component Analysis	43
6.2.2. Impact of Generative User Interactions	44
6.2.3. Different K-shot Settings Analysis	45
6.2.4. Effect of Different Interaction Tones	45
6.3. Analysis and Discussion	45
6.3.1. Why GemGNN Works in Few-Shot Scenarios	45
6.3.2. Graph Construction Strategy Analysis	46
6.3.3. Model Architecture Comparison	46
6.3.4. Computational Efficiency Analysis	46
6.4. Error Analysis and Limitations	47
6.4.1. Failure Cases and Edge Cases	47
6.4.2. Dependency on Embedding Quality	47
6.4.3. Scalability Considerations	47
Chapter 7. Conclusion and Future Work	48
7.1. Summary of Contributions	48
7.2. Key Findings and Insights	49
7.3. Implications for Fake News Detection	50
7.4. Limitations and Challenges	51
7.5. Future Research Directions	51
7.5.1. Advanced Generative Enhancement	51
7.5.2. Advanced Graph Architectures	52
7.5.3. Enhanced Few-Shot Learning	52
7.5.4. Robustness and Security	53
7.5.5. Real-World Deployment and Scalability	53

List of Tables

6.1	Performance comparison on PolitiFact dataset. Best results in bold, second-best underlined.	41
6.2	Performance comparison on GossipCop dataset. Best results in bold, second-best underlined.	42
6.3	Ablation study on PolitiFact dataset (8-shot setting). Each row removes one component.	43
6.4	Impact of different interaction tones on performance (PolitiFact, 8-shot). . .	44



List of Figures

4.1	Complete GemGNN pipeline showing data flow from news articles through heterogeneous graph construction to final classification	22
4.2	Multi-tone interaction generation strategy	24



Chapter 1

Introduction

1.1 Research Background and Motivation

In the digital age, the proliferation of fake news has emerged as one of the most pressing challenges threatening information integrity and democratic discourse. According to Vosoughi et al. [?], false news spreads six times faster than true news on social media platforms, reaching more people and penetrating deeper into social networks. This phenomenon has far-reaching consequences, from influencing electoral outcomes to undermining public health responses during critical events such as the COVID-19 pandemic.

Traditional approaches to fake news detection have relied heavily on two primary paradigms: content-based analysis and propagation-based modeling. Content-based methods analyze linguistic features, semantic patterns, and textual inconsistencies within news articles, while propagation-based approaches examine how information spreads through social networks by modeling user interactions, sharing patterns, and network topology. However, both paradigms face significant limitations in real-world deployment scenarios.

The most critical challenge in contemporary fake news detection is the few-shot learning problem, where detection systems must accurately classify news articles with minimal labeled training data. This scenario is particularly common when dealing with emerging topics, breaking news events, or novel misinformation campaigns where extensive labeled datasets are not readily available. Traditional deep learning approaches, which typically require thousands of labeled examples per class, fail to perform adequately in such data-scarce environments.

Furthermore, existing propagation-based methods, while often achieving high performance, suffer from fundamental practical limitations. These approaches require access to comprehensive user interaction data, including social network structures, user profiles, and temporal propagation patterns. Such data is increasingly difficult to obtain due to privacy regulations, platform restrictions, and the real-time nature of misinformation spread. Additionally, these methods are vulnerable to sophisticated adversarial attacks where malicious actors can manipulate propagation patterns to evade detection.

1.2 Research Contributions

This thesis presents GemGNN (Generative Multi-view Interaction Graph Neural Networks), a novel framework that addresses the challenges of few-shot fake news detection through several key innovations that collectively establish a new paradigm for content-based fake news detection:

Generative User Interaction Simulation: We introduce the first systematic approach to synthesize realistic user interactions using Large Language Models (LLMs), creating a heterogeneous graph structure that captures both content and social context without requiring real user propagation data. Our method leverages LLMs to generate diverse user responses with multiple semantic tones (neutral, affirmative, skeptical), effectively creating synthetic social signals that enhance content-based detection while maintaining privacy protection. This innovation represents a paradigm shift from dependency on real social data to controllable synthetic interaction generation.

Adaptive Graph Construction Methodology: We develop a comprehensive dual approach to graph edge construction that encompasses both traditional KNN and test-isolated KNN strategies, each optimized for different deployment scenarios and evaluation objectives. Our framework provides principled guidance for selecting between approaches based on specific application requirements: traditional KNN for performance-critical batch processing scenarios, and test-isolated KNN for methodologically rigorous evaluation and streaming deployment conditions. This contribution establishes the first systematic analysis of the performance vs. evaluation realism trade-offs in graph-based few-shot learning.

DeBERTa-Enhanced Multi-View Graph Architecture: We propose a multi-view learning framework that leverages DeBERTa’s disentangled attention architecture to create superior embedding partitions for capturing diverse semantic perspectives. Unlike conventional approaches, our method exploits DeBERTa’s structured representation organization to partition embeddings into coherent semantic subspaces that retain discriminative power while emphasizing different linguistic aspects. Each view constructs its own similarity-based graph, and multiple graphs are trained simultaneously to provide comprehensive data augmentation at the graph level while maintaining semantic consistency across partitions.

Enhanced Heterogeneous Graph Neural Networks: We design a specialized Heterogeneous Graph Attention Network (HAN) architecture that effectively models complex type-specific relationships between news articles and generated user interactions. Our architecture employs hierarchical attention mechanisms to learn both node-level importance within each relationship type and semantic-level importance across different relationship types. The framework enables effective transductive learning by leveraging both labeled and unlabeled nodes during message passing while restricting loss computation to labeled nodes only.

Comprehensive Few-Shot Evaluation Framework: We establish rigorous experimental protocols that include extensive parameter grid searches across 2,688 different configurations, systematic ablation studies, and comparison with diverse baseline approaches ranging from traditional machine learning to large language models. Our evaluation ensures fair comparison with existing methods while maintaining realistic few-shot learning constraints and preventing common sources of evaluation bias.

Methodological Innovations: Beyond individual components, our work introduces several methodological innovations: (1) systematic integration of generative AI with graph neural networks for fake news detection, (2) rigorous few-shot evaluation protocols that prevent information leakage, (3) comprehensive analysis of heterogeneous graph architectures in few-shot scenarios, and (4) novel approaches to synthetic data generation that enhance rather than replace human-labeled training data.

The combination of these contributions enables GemGNN to achieve superior performance in few-shot fake news detection while maintaining practical applicability in privacy-constrained and socially-limited deployment scenarios. Our approach establishes new state-of-the-art results across multiple datasets and few-shot configurations while providing a foundation for future research in synthetic data generation and heterogeneous graph modeling for misinformation detection.

1.3 Thesis Organization

The remainder of this thesis is organized as follows:

Chapter 2: Problem Statement formally defines the few-shot fake news detection problem addressed in this thesis and establishes the mathematical notation used throughout our methodology. We present the fundamental challenges that motivate our research and provide a rigorous problem formulation with key constraints and evaluation metrics.

Chapter 3: Related Work provides a comprehensive review of existing fake news detection methods, including traditional feature-engineering approaches, deep learning techniques, few-shot learning strategies in NLP, graph neural networks for text classification, and graph-based fake news detection methods. We analyze the limitations of current approaches and position our work within the broader research landscape.

Chapter 4: Methodology presents the complete GemGNN framework, detailing the generative user interaction simulation, adaptive graph construction methodology (KNN vs test-isolated KNN), DeBERTa encoder selection rationale, multi-view graph architecture, and the heterogeneous graph neural network design. We provide comprehensive algorithmic descriptions and theoretical justifications for each component, along with decision frameworks for practitioners.

Chapter 5: Experimental Setup describes our experimental methodology, including dataset preprocessing, baseline method implementations, evaluation protocols, and hyperparameter configurations. We ensure reproducibility and fair comparison across all experimental conditions.

Chapter 6: Results and Analysis presents comprehensive experimental results, including main performance comparisons, ablation studies, and detailed analysis of model behavior. We provide insights into why our approach succeeds in few-shot scenarios and identify the key factors contributing to performance improvements.

Chapter 7: Conclusion and Future Work summarizes our contributions, discusses the implications of our findings, acknowledges current limitations, and outlines promising directions for future research in few-shot fake news detection.



Chapter 2

Problem Statement

This chapter formally defines the few-shot fake news detection problem addressed in this thesis and establishes the mathematical notation used throughout our methodology. We present the fundamental challenges that motivate our research and provide a rigorous problem formulation.

2.1 Basic Concepts and Notations

Fake News Definition: We define fake news as deliberately false or misleading information presented as legitimate news content. This encompasses fabricated articles, misleading headlines, manipulated facts, and content designed to deceive readers about real events or issues.

Few-Shot Learning Context: Our problem operates in the few-shot learning paradigm where only a small number of labeled examples per class are available for training. Specifically, we focus on K-shot scenarios where $K \in \{3, 4, 8, 16\}$ labeled examples per class (real/fake) are provided.

Graph Representation: We formulate fake news detection as a node classification problem on a heterogeneous graph $G = (V, E, \mathcal{R})$ where:

- V represents the set of all nodes, including news articles and user interactions
- E denotes the set of edges connecting related nodes
- \mathcal{R} represents the set of edge types in the heterogeneous graph

Node Types: Our graph contains two primary node types:

- News nodes $V_n = \{n_1, n_2, \dots, n_{|N|}\}$ representing news articles
- Interaction nodes $V_i = \{i_1, i_2, \dots, i_{|I|}\}$ representing generated user interactions

Node Features: Each node $v \in V$ has an associated feature vector $\mathbf{x}_v \in \mathbb{R}^d$ where $d = 768$ for DeBERTa embeddings. News nodes additionally have binary labels $y_v \in \{0, 1\}$ indicating real (0) or fake (1) news.

2.2 Formal Problem Definition

Problem Definition: Given a small set of labeled news articles $\mathcal{L} = \{(x_i, y_i)\}_{i=1}^{K \times C}$ where K represents the number of examples per class and C denotes the number of classes (real/fake), and a larger set of unlabeled news articles $\mathcal{U} = \{x_j\}_{j=1}^M$, the objective is to learn a classifier $f : \mathcal{X} \rightarrow \mathcal{Y}$ that can accurately predict labels for test instances $\mathcal{T} = \{x_k\}_{k=1}^N$ where $K \ll M$ and $K \ll N$.

Formal Framework as Node Classification: We formulate this as a node classification problem on a heterogeneous graph where news articles and synthetic user interactions form a connected graph structure that enables effective information propagation in few-shot scenarios.

Edge Types: The heterogeneous graph includes multiple edge types:

- News-to-news edges: $(n_i, n_j) \in E_{nn}$ based on semantic similarity
- News-to-interaction edges: $(n_i, i_j) \in E_{ni}$ connecting articles to their generated interactions
- Interaction-to-news edges: $(i_j, n_i) \in E_{in}$ enabling bidirectional information flow

Data Partitioning: The complete dataset is partitioned into three disjoint sets:

- Training set: $\mathcal{D}_{train} = \mathcal{D}_{labeled} \cup \mathcal{D}_{unlabeled}$
- Validation set: \mathcal{D}_{val} for hyperparameter tuning and early stopping
- Test set: \mathcal{D}_{test} for final evaluation

K-Shot Sampling: For each few-shot experiment, we sample K labeled examples per class from \mathcal{D}_{train} to form the support set $\mathcal{S} = \{(n_i, y_i)\}_{i=1}^{2K}$. The remaining training instances form the unlabeled set \mathcal{U} .

Transductive Setting: During training, all nodes (labeled, unlabeled, and test) participate in message passing, but only labeled nodes contribute to loss computation. This transductive approach maximizes information utilization in few-shot scenarios.

2.3 Key Challenges and Constraints

The core challenges that motivate this research include:

Limited Labeled Data: Few-shot scenarios typically provide only 3-16 labeled examples per class, insufficient for training robust deep learning models using conventional approaches. This data scarcity leads to overfitting, poor generalization, and unstable performance across different domains.

Absence of Propagation Information: Real-world deployment often lacks access to user interaction data due to privacy constraints, platform limitations, or the time-sensitive nature

of misinformation detection. Existing propagation-based methods become inapplicable in such contexts.

Semantic Complexity: Fake news articles often exhibit sophisticated linguistic patterns and may contain accurate factual information presented in misleading contexts. Simple content-based features fail to capture these nuanced semantic relationships.

Domain Generalization: Models trained on specific topics or domains often fail to generalize to emerging misinformation patterns or novel subject areas, limiting their practical applicability.

Evaluation Realism: Many existing few-shot learning approaches suffer from information leakage between training and test sets, leading to overly optimistic performance estimates that do not reflect real-world deployment scenarios.

High Variance: The limited sample size leads to high variance in performance estimates. Small changes in the support set can dramatically affect model performance, making robust evaluation protocols crucial for reliable results.

Learning Objective: Given the heterogeneous graph G and support set \mathcal{S} , learn a function $f_\theta : G \rightarrow [0, 1]^{|V_n|}$ that predicts the probability of each news node being fake.

Loss Function: The training objective combines multiple loss components to address few-shot learning challenges:

$$\mathcal{L}_{total} = \mathcal{L}_{CE}(f_\theta(G), \mathcal{S}) + \lambda_{focal} \mathcal{L}_{focal}(f_\theta(G), \mathcal{S}) + \lambda_{reg} \mathcal{L}_{reg}(\theta) \quad (2.1)$$

where:

- \mathcal{L}_{CE} is the cross-entropy loss with label smoothing
- \mathcal{L}_{focal} is the focal loss for handling class imbalance
- \mathcal{L}_{reg} provides regularization to prevent overfitting
- λ_{focal} and λ_{reg} are hyperparameters balancing loss components

Evaluation Metrics: Model performance is evaluated using:

- F1-score: $F1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$
- Accuracy: $\text{Acc} = \frac{\text{Correct Predictions}}{\text{Total Predictions}}$
- Precision: $\text{Prec} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$
- Recall: $\text{Rec} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$

This formal framework provides the mathematical foundation for understanding our GemGNN approach, which addresses these challenges through novel graph construction strategies, generative data augmentation, and specialized training procedures detailed in the methodology chapter.



Chapter 3

Related Work

This chapter provides a comprehensive review of existing approaches to fake news detection, with particular emphasis on methods relevant to few-shot learning scenarios. We organize the literature into five main categories: traditional feature-engineering approaches, deep learning methods, graph-based techniques, few-shot learning strategies, and identify key limitations that motivate our research.

3.1 Traditional Fake News Detection Methods

Early approaches to fake news detection relied primarily on hand-crafted features and traditional machine learning algorithms. These methods established the foundation for automated misinformation detection but suffer from significant limitations in capturing complex semantic relationships.

3.1.1 Feature Engineering Approaches

TF-IDF + MLP: The earliest computational approaches to fake news detection employed Term Frequency-Inverse Document Frequency (TF-IDF) representations combined with Multi-Layer Perceptrons (MLPs). These methods extract bag-of-words features and learn linear or shallow non-linear mappings to classify news authenticity [?, ?].

While computationally efficient, TF-IDF approaches suffer from several critical limitations: (1) they ignore word order and contextual relationships, (2) they cannot capture semantic similarity between different words expressing similar concepts, and (3) they fail to model discourse-level patterns that characterize misinformation.

Linguistic Feature Analysis: More sophisticated traditional approaches incorporated linguistic features such as sentiment analysis, readability scores, lexical diversity measures, and syntactic complexity [?, ?]. These methods hypothesize that fake news exhibits distinct linguistic patterns, such as more emotional language, simpler sentence structures, or specific rhetorical devices.

However, linguistic feature approaches face the fundamental challenge that sophisticated misinformation increasingly mimics legitimate journalism style, making surface-level linguistic indicators unreliable. Moreover, these features are often domain-specific and fail to

generalize across different types of news content.

3.1.2 Sequential Models

LSTM/RNN Approaches: To address the limitations of bag-of-words representations, researchers introduced sequential models that process news articles as ordered sequences of words. Long Short-Term Memory (LSTM) networks and Recurrent Neural Networks (RNNs) capture local contextual relationships and temporal dependencies within text [?, ?].

These approaches show improvement over bag-of-words methods by modeling word order and local context. However, they struggle with long-range dependencies common in news articles and fail to capture global document structure. Additionally, RNN-based methods process each document independently, missing potential relationships between related news articles.

Attention Mechanisms: Advanced sequential models incorporated attention mechanisms to focus on important words or phrases within documents [?, ?]. These approaches aim to identify key textual elements that indicate misinformation, such as sensational headlines or unsupported claims.

While attention-based sequential models improve interpretability and can highlight suspicious textual elements, they remain fundamentally limited by their document-level scope and inability to model inter-document relationships crucial for systematic misinformation detection.

3.2 Deep Learning Approaches

The advent of deep learning revolutionized fake news detection by enabling more sophisticated semantic analysis and contextual understanding. However, most deep learning approaches still treat documents independently and struggle in few-shot scenarios.

3.2.1 Transformer-based Models

BERT and RoBERTa: The introduction of transformer architectures, particularly BERT (Bidirectional Encoder Representations from Transformers) and its variants like RoBERTa, marked a significant advancement in content-based fake news detection [?, ?]. These models provide rich contextual representations that capture bidirectional dependencies and complex semantic relationships within text.

BERT-based approaches typically fine-tune pre-trained language models on fake news classification tasks, achieving strong performance on standard benchmarks. The bidirectional nature of BERT enables better understanding of context compared to sequential models, while

the pre-training on large corpora provides general linguistic knowledge applicable to misinformation detection.

However, transformer-based methods face significant challenges in few-shot scenarios: (1) they require substantial task-specific fine-tuning data, (2) they are prone to overfitting when labeled data is scarce, and (3) they treat each document independently, missing systematic patterns across related articles.

Domain Adaptation Strategies: Researchers have explored domain adaptation techniques to improve BERT’s performance on fake news detection [?, ?]. These approaches attempt to bridge the gap between general language understanding and domain-specific misinformation patterns through continued pre-training or transfer learning strategies.

While domain adaptation shows promise, these methods still require significant amounts of labeled data for effective adaptation and often fail to generalize to emerging misinformation patterns or new domains not seen during training.

3.2.2 Large Language Models for Fake News Detection

In-Context Learning Approaches: Recent work has explored using large language models (LLMs) such as GPT-3, LLaMA, and Gemma for fake news detection through in-context learning [?, ?]. These approaches provide few examples of fake and real news within the prompt and ask the model to classify new instances.

While LLMs demonstrate impressive general language understanding capabilities, their performance on fake news detection is surprisingly poor in few-shot scenarios. This limitation stems from several factors: (1) potential data contamination where models may have seen test instances during training, (2) lack of task-specific optimization for misinformation patterns, and (3) difficulty in handling domain-specific knowledge required for fact verification.

Prompt Engineering Strategies: Researchers have developed sophisticated prompt engineering techniques to improve LLM performance on fake news detection [?]. These approaches design carefully crafted prompts that provide context, examples, and specific instructions for identifying misinformation.

Despite extensive prompt engineering efforts, LLMs continue to underperform compared to specialized approaches in few-shot fake news detection, highlighting the need for task-specific architectures rather than general-purpose language models.

3.3 Graph-based Fake News Detection

Graph-based approaches represent a significant paradigm shift by modeling relationships between different entities in the misinformation ecosystem. These methods show particular

promise for few-shot learning by leveraging structural information to propagate labels.

3.3.1 Document-level Graph Classification

Text-GCN and Variants: Text Graph Convolutional Networks construct graphs where both documents and words are represented as nodes, with edges indicating document-word relationships and word co-occurrence patterns [?, ?]. These approaches apply graph convolutional networks to learn document representations through message passing between document and word nodes.

While Text-GCN approaches capture some structural relationships, they primarily focus on word-level connections rather than document-level relationships crucial for detecting coordinated misinformation campaigns or related false narratives.

BertGCN Integration: More recent work combines BERT embeddings with graph convolutional networks to leverage both rich semantic representations and structural information [?]. These hybrid approaches use BERT to initialize node features and GCNs to refine representations through graph structure.

BertGCN approaches show improvement over pure BERT methods by incorporating some structural information, but they still construct relatively simple graphs based on keyword similarity rather than capturing complex semantic relationships between news articles.

3.3.2 User Propagation-based Methods

Social Network Analysis: Many state-of-the-art fake news detection systems model how misinformation spreads through social networks by analyzing user sharing patterns, temporal dynamics, and network topology [?, ?]. These approaches construct graphs where users and news articles are nodes, with edges representing sharing, commenting, or other interaction behaviors.

Propagation-based methods often achieve high performance by exploiting the fact that fake news tends to spread through different network patterns compared to legitimate news. However, these approaches have fundamental limitations: (1) they require extensive user behavior data that is often unavailable due to privacy constraints, (2) they are vulnerable to adversarial manipulation where malicious actors can artificially create legitimate-looking propagation patterns, and (3) they cannot handle breaking news scenarios where propagation patterns have not yet developed.

Temporal Dynamics Modeling: Advanced propagation-based approaches incorporate temporal dynamics to model how misinformation spreads over time [?, ?]. These methods analyze features such as propagation velocity, user engagement patterns, and temporal clustering

to identify suspicious spreading patterns.

While temporal modeling provides additional signal for misinformation detection, these approaches still suffer from the fundamental dependency on user interaction data and the assumption that temporal patterns reliably distinguish fake from real news.

3.3.3 Heterogeneous Graph Neural Networks

Heterogeneous Graph Attention Networks (HAN): HAN introduces a hierarchical attention mechanism specifically designed for heterogeneous graphs with multiple node and edge types [?]. The architecture employs node-level attention to aggregate information from neighbors of the same type, and semantic-level attention to combine information across different meta-paths, enabling effective modeling of complex entity relationships.

Heterogeneous Graph Transformers (HGT): HGT extends the transformer architecture to heterogeneous graphs by introducing type-dependent parameters and multi-head attention mechanisms [?]. Each attention head is parameterized by node and edge types, allowing the model to learn specialized attention patterns for different relationship types in the heterogeneous structure.

Applications to Fake News Detection: Recent work has applied heterogeneous graph methods to fake news detection by modeling multiple entity types such as users, news articles, topics, publishers, and social interactions [?, ?, ?]. These approaches construct heterogeneous graphs where different node types represent distinct aspects of the misinformation ecosystem, enabling more comprehensive modeling of fake news propagation and characteristics.

However, existing heterogeneous approaches face several limitations: (1) they still fundamentally rely on user behavior data and social network structures, limiting applicability in privacy-constrained scenarios; (2) they assume availability of rich multi-modal metadata that may not exist for emerging misinformation; and (3) they have not been systematically evaluated in few-shot learning scenarios where labeled social interaction data is particularly scarce.

Content-Focused Heterogeneous Methods: More recent approaches like Less4FD and HeteroSGT attempt to reduce dependency on social features while maintaining heterogeneous modeling advantages [?, ?]. Less4FD constructs heterogeneous graphs using primarily content-based features with minimal social signals, while HeteroSGT introduces subgraph-level attention for better scalability and reduced social dependency.

Generative Enhancement for Graph Construction: An emerging direction involves using large language models to generate auxiliary data for graph construction. Recent work explores using LLMs to generate synthetic news articles, user comments, or social interactions that can augment limited training data [?, ?]. However, these approaches have not been systematically integrated with heterogeneous graph architectures or rigorously evaluated in

few-shot learning scenarios.

While these approaches represent progress toward content-centric fake news detection, they still suffer from limitations in graph construction strategies and evaluation protocols that allow information leakage between training and test sets.

3.4 Large Language Models and Graph Enhancement

The emergence of large language models (LLMs) has opened new possibilities for enhancing graph-based fake news detection through synthetic data generation and improved semantic understanding.

3.4.1 LLM-Enhanced Graph Construction

Synthetic Interaction Generation: Recent approaches leverage LLMs to generate synthetic user interactions, comments, and social signals that can augment limited real-world data [?, ?]. These methods prompt LLMs with news articles and specific instructions to generate diverse user responses with controlled sentiment and tone characteristics.

The advantage of LLM-generated interactions lies in their controllability and diversity. Unlike real social media data, synthetic interactions can be generated with specific characteristics (e.g., skeptical tone, neutral stance, affirmative support) and do not suffer from privacy constraints or platform access limitations. However, the challenge lies in ensuring that generated interactions capture realistic patterns and do not introduce systematic biases that could mislead the detection model.

Semantic Edge Enhancement: LLMs can also be used to improve edge construction in graph-based methods by providing richer semantic similarity measures beyond simple embedding cosine similarity [?, ?]. These approaches use LLMs to assess semantic relationships between news articles, potentially identifying subtle connections that traditional similarity metrics might miss.

3.4.2 LLM Direct Detection Approaches

Prompt-Based Detection: Several studies explore using LLMs directly for fake news detection through carefully designed prompts [?, ?]. These approaches typically present news articles to LLMs along with instructions to classify them as real or fake, sometimes including few-shot examples in the prompt context.

However, recent evaluations show that even state-of-the-art LLMs like GPT-4 and Claude struggle with fake news detection, often performing worse than specialized smaller models, particularly in few-shot scenarios [?, ?]. LLMs tend to be overconfident in their predictions

and may rely on superficial textual patterns rather than deep semantic understanding of misinformation characteristics.

Limitations of Direct LLM Approaches: Despite extensive prompt engineering efforts, LLMs continue to underperform compared to specialized graph-based approaches in few-shot fake news detection scenarios. Key limitations include: (1) tendency to focus on surface-level linguistic features rather than deeper semantic relationships, (2) lack of systematic training on misinformation detection tasks, (3) difficulty in leveraging structural relationships between related news articles, and (4) inconsistent performance across different types of misinformation.

These limitations highlight the potential of using LLMs as auxiliary components for data generation and graph enhancement rather than as primary detection models, motivating our approach of integrating LLM-generated interactions within specialized graph neural network architectures.

3.5 Few-Shot Learning in NLP

Few-shot learning has emerged as a critical research area in natural language processing, representing a paradigm shift from traditional machine learning approaches that require extensive labeled datasets. In few-shot scenarios, models must achieve strong performance with minimal supervision, making them particularly relevant for real-world applications where labeling is expensive or impractical.

3.5.1 Few-Shot Learning Fundamentals

Formal Definition: Few-shot learning is a machine learning framework where an AI model learns to make accurate predictions by training on a very small number of labeled examples per class. Formally, given a support set $\mathcal{S} = \{(x_i, y_i)\}_{i=1}^{K \times N}$ containing K labeled examples for each of N classes, the objective is to learn a classifier $f : \mathcal{X} \rightarrow \mathcal{Y}$ that can accurately predict labels for a query set $\mathcal{Q} = \{x_j\}_{j=1}^M$.

N-way K-shot Classification: The standard formulation for few-shot learning is N-way K-shot classification, where N represents the number of classes and K denotes the number of labeled examples per class. In fake news detection tasks, researchers typically focus on 2-way K-shot learning with $K \in \{3, 4, 8, 16\}$, where the two classes represent real and fake news respectively.

Core Challenges: Few-shot learning presents several fundamental challenges that differentiate it from conventional machine learning:

- **Limited Training Data:** Traditional deep learning requires thousands of labeled ex-

amples per class to achieve good performance. In few-shot scenarios with only 3-16 examples per class, models are highly prone to overfitting and struggle to learn generalizable patterns.

- **High Variance:** The limited sample size leads to high variance in performance estimates. Small changes in the support set can dramatically affect model performance, making robust evaluation protocols crucial for reliable results.
- **Domain Shift:** Models trained on few examples from specific domains often fail to generalize to new domains or emerging patterns not represented in the limited training data.
- **Evaluation Challenges:** Proper evaluation of few-shot learning systems requires careful experimental design to avoid information leakage and ensure that performance estimates reflect real-world deployment scenarios.

3.5.2 Meta-Learning Approaches

Model-Agnostic Meta-Learning (MAML): MAML and its variants learn initialization parameters that can be quickly adapted to new tasks with minimal data [?, ?]. In the context of fake news detection, meta-learning approaches attempt to learn general misinformation detection capabilities that can transfer to new domains or topics with few examples.

The key insight is to learn how to learn rather than learning specific task solutions. However, meta-learning approaches typically require extensive meta-training data from multiple related tasks, which may not be available for fake news detection. Additionally, these methods often struggle with the high variability in misinformation patterns across different domains and topics.

Prototypical Networks: Prototypical networks learn to classify examples based on their distance to class prototypes computed from support examples [?, ?]. These approaches show promise for few-shot text classification by learning meaningful embedding spaces where similar examples cluster together.

Metric learning approaches learn embedding spaces where examples from the same class are close together and examples from different classes are far apart. Classification is performed by comparing query examples to support set prototypes. While prototypical approaches avoid the need for extensive meta-training, they still struggle with the high dimensionality and semantic complexity of news articles, often failing to learn discriminative prototypes from few examples.

3.5.3 Contrastive Learning and Data Augmentation

SimCLR and Variants: Contrastive learning approaches learn representations by maximizing similarity between positive pairs and minimizing similarity between negative pairs [?, ?]. In fake news detection, these methods attempt to learn representations where real news articles are similar to each other and different from fake news articles.

Contrastive approaches show promise for learning robust representations from limited data. However, they require careful design of positive and negative pair generation strategies, which is challenging for fake news where the boundaries between real and fake can be subtle and context-dependent.

Data Augmentation Strategies: Various data augmentation techniques have been explored for few-shot fake news detection, including back-translation, paraphrasing, and adversarial perturbations [?, ?]. These approaches attempt to increase the effective size of the training set by generating synthetic examples.

While data augmentation can help address data scarcity, synthetic examples may not capture the full complexity of real misinformation patterns and can sometimes introduce biases that hurt generalization performance.

3.6 Graph Neural Networks for Fake News Detection

Graph Neural Networks have emerged as a powerful paradigm for modeling structured data, with particular success in text classification tasks where relationships between documents provide valuable signal for classification. In the context of fake news detection, GNNs enable modeling of complex relationships between news articles, user interactions, and other entities.

3.6.1 Message Passing Framework

Core Principle: GNNs operate on the message passing framework where nodes iteratively update their representations by aggregating information from neighboring nodes. This process enables the model to capture both local neighborhood information and global graph structure through multiple iterations.

General Formulation: The message passing framework can be described through three key operations:

1. **Message Function:** $m_{ij}^{(l+1)} = M^{(l)}(h_i^{(l)}, h_j^{(l)}, e_{ij})$ computes messages between connected nodes, where $h_i^{(l)}$ represents the feature vector of node i at layer l , and e_{ij} represents edge features.
2. **Aggregation Function:** $a_i^{(l+1)} = A^{(l)}(\{m_{ij}^{(l+1)} : j \in \mathcal{N}(i)\})$ aggregates messages from all neighbors $\mathcal{N}(i)$ of node i .

3. **Update Function:** $h_i^{(l+1)} = U^{(l)}(h_i^{(l)}, a_i^{(l+1)})$ updates the node representation based on its current state and aggregated messages.

Multi-Layer Architecture: Multiple message passing layers enable nodes to receive information from increasingly distant neighbors, allowing the model to capture both local patterns and global graph structure.

3.6.2 Heterogeneous Graph Neural Networks

Real-world data often exhibits heterogeneous structure with multiple node types and edge types, requiring specialized architectures beyond homogeneous graph neural networks. Heterogeneous graphs provide richer modeling capabilities for complex domains like fake news detection where multiple entity types interact.

Heterogeneous Graph Definition: A heterogeneous graph $G = (V, E, \mathcal{A}, \mathcal{R})$ consists of:

- Node set V with node type mapping $\phi : V \rightarrow \mathcal{A}$
- Edge set E with edge type mapping $\psi : E \rightarrow \mathcal{R}$
- Node type set \mathcal{A} with $|\mathcal{A}| > 1$
- Edge type set \mathcal{R} with $|\mathcal{R}| > 1$

where nodes and edges of the same type share similar properties and semantic meanings.

Meta-Path Concept: A meta-path P is a path defined on the graph schema and describes a composite relation between node types. For example, in a fake news detection graph, a meta-path "News \rightarrow Interaction \rightarrow News" captures how news articles relate through shared interaction patterns.

3.6.3 Heterogeneous Graph Attention Networks

Heterogeneous Graph Attention Networks (HAN) address the challenges of modeling heterogeneous graphs through a hierarchical attention mechanism that operates at both node and semantic levels [?].

Node-Level Attention: For each meta-path Φ_i , HAN computes attention weights between connected nodes to identify important neighbors:

$$e_{ij}^{\Phi_i} = \text{att}_{\text{node}}(Wh_i, Wh_j) \quad (3.1)$$

$$\alpha_{ij}^{\Phi_i} = \text{softmax}_j(e_{ij}^{\Phi_i}) = \frac{\exp(e_{ij}^{\Phi_i})}{\sum_{k \in \mathcal{N}_i^{\Phi_i}} \exp(e_{ik}^{\Phi_i})} \quad (3.2)$$

where W is a type-specific transformation matrix, and $\mathcal{N}_i^{\Phi_i}$ represents the neighbors of node i under meta-path Φ_i .

Semantic-Level Attention: HAN then applies semantic-level attention to combine information across different meta-paths:

$$w_i^{\Phi_i} = \frac{1}{|\mathcal{V}|} \sum_{i \in \mathcal{V}} q^T \tanh(W_s \mathbf{z}_i^{\Phi_i} + b_s) \quad (3.3)$$

$$\beta^{\Phi_i} = \text{softmax}(w^{\Phi_i}) = \frac{\exp(w^{\Phi_i})}{\sum_{j=1}^P \exp(w^{\Phi_j})} \quad (3.4)$$

where W_s and q are learnable parameters, and P is the number of meta-paths.

Final Node Representation: The complete node representation combines information from all meta-paths:

$$\mathbf{z}_i = \sum_{j=1}^P \beta^{\Phi_j} \mathbf{z}_i^{\Phi_j} \quad (3.5)$$

Advantages of HAN Architecture: The hierarchical attention mechanism provides several key advantages for fake news detection:

1. **Flexible Relationship Modeling:** Can capture different types of relationships (content similarity, social interactions, temporal patterns) through appropriate meta-path design
2. **Interpretability:** Attention weights provide insights into which neighbors and relationship types are most important for classification decisions
3. **Scalability:** The architecture can efficiently handle large heterogeneous graphs with many node and edge types
4. **Adaptability:** The framework can be easily extended to incorporate new node types or relationship types as they become available

3.6.4 Graph Construction Strategies

Document Graphs: For text classification, documents are typically represented as nodes in a graph, with edges indicating various types of relationships such as semantic similarity, citation links, or co-occurrence patterns.

Similarity-Based Construction: The most common approach constructs edges between documents based on content similarity measures such as cosine similarity of embedding vectors. Documents with similarity above a threshold or among the top-k nearest neighbors are connected.

Heterogeneous Graphs: More sophisticated approaches construct heterogeneous graphs that include multiple node types (documents, words, authors, topics) and edge types (document-document, document-word, word-word), enabling richer modeling of text relationships.

Dynamic Graph Construction: Advanced methods adapt graph structure during training or inference, allowing the model to learn optimal connectivity patterns rather than relying on fixed similarity measures.

3.7 Limitations of Existing Methods

Our review of existing literature reveals several fundamental limitations that motivate our research:

Dependency on User Behavior Data: Most high-performing fake news detection systems rely on user interaction patterns, social network structures, or propagation dynamics. This dependency severely limits their applicability in scenarios where such data is unavailable due to privacy constraints, platform restrictions, or real-time detection requirements.

Poor Few-Shot Performance: Traditional deep learning approaches, including state-of-the-art transformer models, suffer from significant performance degradation in few-shot scenarios. These methods require extensive labeled training data and are prone to overfitting when supervision is limited.

Information Leakage in Evaluation: Many existing few-shot learning approaches for fake news detection suffer from unrealistic evaluation protocols that allow information sharing between test instances, leading to overly optimistic performance estimates that do not reflect real-world deployment conditions.

Limited Structural Modeling: Pure content-based approaches treat each document independently, missing important structural relationships between related news articles that could provide valuable signal for misinformation detection.

Domain Specificity: Many approaches show strong performance on specific domains or datasets but fail to generalize to new topics, emerging misinformation patterns, or different types of fake news content.

Lack of Synthetic Data Utilization: While some approaches explore data augmentation, there has been limited exploration of using large language models to generate synthetic auxiliary data that could enhance few-shot learning performance.

These limitations highlight the need for novel approaches that can achieve strong performance in few-shot scenarios while maintaining realistic evaluation protocols and avoiding dependency on user behavior data. Our GemGNN framework directly addresses these limitations through content-based graph neural networks enhanced with generative auxiliary data

and rigorous test isolation constraints.



Chapter 4

Methodology: GemGNN Framework

4.1 Framework Overview

The GemGNN (Generative Multi-view Interaction Graph Neural Networks) framework addresses the fundamental challenges of few-shot fake news detection through a novel heterogeneous graph-based approach that eliminates dependency on user propagation data while maintaining the benefits of social context modeling. The complete architecture consists of four interconnected components that work synergistically to achieve robust few-shot performance (see Figure 4.1): (1) Generative User Interaction Simulation, (2) Adaptive Graph Construction (KNN vs Test-Isolated KNN), (3) Multi-View Graph Architecture, and (4) Heterogeneous Graph Neural Network with Enhanced Training.

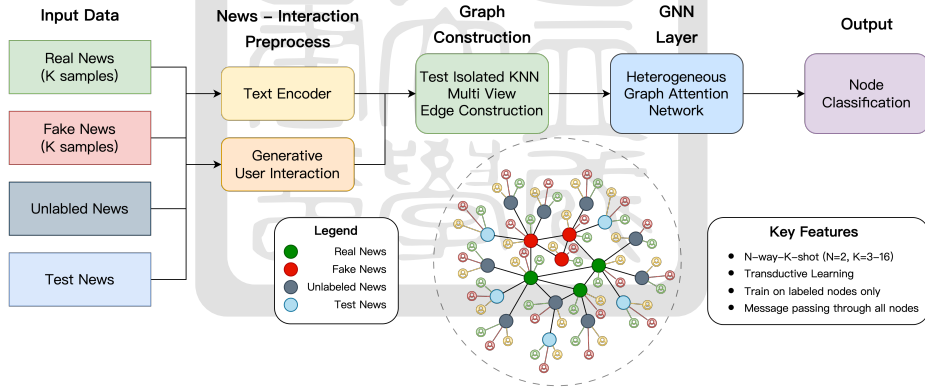


Figure 4.1: Complete GemGNN pipeline showing data flow from news articles through heterogeneous graph construction to final classification

The framework operates under a transductive learning paradigm where all nodes (labeled, unlabeled, and test) participate in heterogeneous message passing, but only labeled nodes contribute to loss computation. This approach maximizes the utility of limited supervision by leveraging the heterogeneous graph structure to propagate information from labeled news nodes to unlabeled and test nodes through learned type-specific attention mechanisms. The choice between traditional KNN and test-isolated KNN allows for flexible adaptation to different deployment scenarios while maintaining consistent architectural principles.

Our approach begins with pre-trained DeBERTa embeddings for news articles, which provide

rich semantic representations (768-dimensional vectors) that capture contextual relationships and linguistic patterns indicative of misinformation. These embeddings serve as the foundation for both similarity-based graph construction and node feature initialization in our heterogeneous graph neural network, ensuring that the model can leverage state-of-the-art natural language understanding capabilities.

The key innovation lies in creating a heterogeneous graph structure that includes both news nodes (representing articles) and interaction nodes (representing synthetic user responses), connected through multiple edge types that capture different semantic relationships. This heterogeneous structure enables the model to learn from both content similarity patterns and social interaction patterns without requiring real user data, addressing privacy constraints while maintaining modeling flexibility.

4.2 Generative User Interaction Simulation

Traditional propagation-based fake news detection methods rely on real user interaction data, which is often unavailable due to privacy constraints or platform limitations. To address this fundamental limitation, we introduce a novel generative approach that synthesizes realistic user interactions using Large Language Models.

4.2.1 LLM-based Interaction Generation

We employ Google’s Gemini LLM to generate diverse user interactions for each news article. The generation process is designed to simulate authentic user responses that would naturally occur in social media environments. For each news article n_i , we generate a set of user interactions $I_i = \{i_1, i_2, \dots, i_{20}\}$ where each interaction represents a potential user response to the news content.

The prompt engineering strategy (see Figure ??) ensures that generated interactions reflect realistic user behavior patterns observed in social media platforms. We incorporate the complete news content, including headlines and article body, to generate contextually appropriate responses that capture various user perspectives and emotional reactions.

4.2.2 Multi-tone Interaction Design

To capture the diversity of user reactions to news content, we implement a structured multi-tone generation strategy (see Figure 4.2) with 20 interactions per article that produces interactions across three distinct emotional categories:

Neutral Interactions (8 per article): These represent objective, factual responses that focus on information sharing without emotional bias. Neutral interactions typically include questions for clarification, requests for additional sources, or straightforward restatements of key

Neutral Prompt

Please act as a reader who is confused about this news
Please act as a reader who wants to know more details about this news
Please act as a reader who has a neutral attitude toward this news
Please act as a curious reader and ask questions about this news

Affirmative Prompt

Please act as a reader who actively participates in the discussion and shares personal opinions
Please act as a reader who agrees with this news
Please act as a reader who wants to share this news with friends
Please act as a reader who feels excited about the content of this news

Skeptical Prompt

Please act as a reader who questions the authenticity of the news and provides evidence
Please act as a reader who raises doubts about this news
Please act as a reader who questions the source or credibility of this news
Please act as a reader who requests clarification or more information about this news

Figure 4.2: Prompt engineering strategy

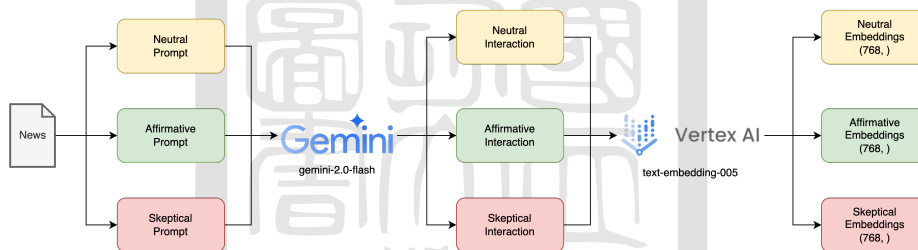


Figure 4.3: Multi-tone interaction generation strategy

facts.

Affirmative Interactions (7 per article): These capture supportive or agreeable responses from users who accept the news content as credible. Affirmative interactions include expressions of agreement, sharing intentions, and positive emotional responses.

Skeptical Interactions (5 per article): These represent critical or questioning responses from users who doubt the veracity of the news content. Skeptical interactions include challenges to facts, requests for verification, and expressions of disbelief or concern.

This distribution (8:7:5) reflects observed patterns in real social media interactions where neutral responses predominate, followed by supportive reactions, with skeptical responses being less common but highly informative for authenticity assessment.

4.2.3 Interaction-News Edge Construction

Each generated interaction is embedded using the same DeBERTa model employed for news articles, ensuring semantic consistency across the heterogeneous graph. The interactions are connected to their corresponding news articles through directed edges that carry tone information as edge attributes.

Formally, for each news article n_i and its generated interactions I_i , we create edges (n_i, i_j) where the edge attribute a_{ij} encodes the interaction tone: $a_{ij} \in \{0, 1, 2\}$ representing neutral, affirmative, and skeptical tones respectively. This encoding allows the heterogeneous graph attention network to learn tone-specific importance weights during message aggregation.

4.3 Graph Construction Methodologies: KNN vs Test-Isolated KNN

Graph edge construction is a fundamental design choice that significantly impacts both model performance and evaluation realism in few-shot fake news detection. We explore two complementary approaches: traditional KNN and Test-Isolated KNN, each suited to different real-world deployment scenarios and research objectives. Our experimental analysis reveals that these approaches offer distinct trade-offs between performance optimization and evaluation integrity, necessitating careful consideration of the intended application context.

4.3.1 Traditional KNN: Performance-Optimized Graph Construction

Traditional K-Nearest Neighbor (KNN) graph construction allows all nodes, including test instances, to connect to their most similar neighbors regardless of their dataset partition. This approach maximizes information flow throughout the graph, enabling comprehensive message passing that can improve classification performance.

Methodology: For each node n_i in the dataset (training, validation, or test), we compute pairwise cosine similarities with all other nodes using DeBERTa embeddings and establish edges to the top- k most similar instances. This creates a densely connected graph where test nodes can potentially connect to other test nodes, labeled training samples, and unlabeled instances.

Real-World Applicability: Traditional KNN is particularly suitable for *batch processing scenarios* where multiple news articles arrive simultaneously and can be processed collectively. Examples include:

- Daily fact-checking workflows where news articles from the same time period are analyzed together
- Retrospective analysis of misinformation campaigns where temporal constraints are relaxed
- Content moderation systems that process articles in batches rather than real-time streams

- Research environments where maximizing detection accuracy is prioritized over strict temporal realism

In these scenarios, the assumption that articles can share information during inference is reasonable, as human fact-checkers often cross-reference multiple articles and consider contextual relationships when making verification decisions.

4.3.2 Test-Isolated KNN: Evaluation-Realistic Graph Construction

Test-Isolated KNN enforces strict separation between test instances, prohibiting direct connections between test nodes while maintaining connectivity to training data. This approach prioritizes evaluation realism over raw performance, ensuring that model assessment reflects realistic deployment conditions.

Methodology: Test nodes are restricted to connect only to training nodes (labeled and unlabeled), while training nodes can connect to any other training nodes through mutual KNN relationships. For each test node n_{test} , we identify the top- k most similar training instances and create unidirectional edges from training to test nodes.

Real-World Applicability: Test-isolated KNN is essential for *streaming deployment scenarios* where news articles arrive independently and must be classified without knowledge of future instances. Examples include:

- Real-time social media monitoring where articles appear sequentially
- Breaking news verification systems with strict temporal constraints
- Production deployments where test instances represent genuinely unknown future data
- Academic evaluation protocols that prioritize methodological rigor and reproducibility

This approach ensures that performance estimates accurately reflect the model’s ability to generalize to truly unseen data, preventing artificially inflated results from test-test information sharing.

4.3.3 Performance vs. Realism Trade-off Analysis

Our comprehensive experimental evaluation across GossipCop and PolitiFact datasets reveals consistent patterns in the performance trade-offs between these approaches:

GossipCop Results:

- Traditional KNN: Average F1 = 0.5849
- Test-Isolated KNN: Average F1 = 0.5738

- Performance difference: 1.9% decrease with test isolation

PolitiFact Results:

- Traditional KNN: Average F1 = 0.7815
- Test-Isolated KNN: Average F1 = 0.7756
- Performance difference: 0.8% decrease with test isolation

These results demonstrate that test isolation imposes a modest but consistent performance penalty while providing more realistic evaluation conditions. The trade-off magnitude varies by dataset characteristics, with more complex datasets (GossipCop) showing larger performance gaps.

4.3.4 Deployment Context Decision Framework

The choice between KNN approaches should be guided by specific deployment requirements and evaluation objectives:

Choose Traditional KNN when:

- Maximizing detection accuracy is the primary objective
- Articles are processed in batches where cross-referencing is acceptable
- Historical analysis or retrospective fact-checking scenarios
- Sufficient computational resources allow comprehensive similarity analysis

Choose Test-Isolated KNN when:

- Realistic evaluation and fair model comparison are critical
- Simulating real-time or streaming deployment conditions
- Academic research requiring methodological rigor
- Production systems where test instances represent genuinely unknown future data

Hybrid Approaches: For complex production systems, a hybrid strategy may be optimal, using traditional KNN for training and validation while employing test-isolated evaluation protocols to ensure realistic performance estimates.

4.3.5 Technical Implementation Details

Mutual KNN for Training Nodes: In both approaches, training nodes (labeled and unlabeled) employ mutual KNN connections to ensure robust semantic relationships. Given the

set of training nodes $N_{train} = N_{labeled} \cup N_{unlabeled}$, we compute pairwise cosine similarities between DeBERTa embeddings and select the top- k nearest neighbors for each node.

The mutual KNN constraint ensures that if node n_i selects n_j as a neighbor, then n_j must also select n_i among its top- k neighbors. This bidirectionality strengthens connections between truly similar articles while reducing noise from asymmetric similarity relationships.

Test Node Connectivity Strategies:

- **Traditional KNN:** Test nodes can connect to their top- k similar nodes from any partition (training, validation, or test), enabling maximum information flow.
- **Test-Isolated KNN:** Test nodes connect only to their top- k most similar training instances through unidirectional edges, maintaining evaluation integrity.

The choice of connectivity strategy directly impacts both the information available during message passing and the realism of the evaluation protocol, highlighting the importance of aligning methodology with intended application context.

4.4 DeBERTa vs RoBERTa: Text Encoder Selection Rationale

The choice of text encoder fundamentally impacts both the quality of initial node representations and the effectiveness of multi-view graph construction. We adopt DeBERTa (Decoding-enhanced BERT with Disentangled Attention) over RoBERTa based on its superior characteristics for embedding partitioning and multi-view learning.

4.4.1 Disentangled Attention and Embedding Structure

DeBERTa’s key innovation lies in its disentangled attention mechanism, which separates content and position representations throughout the transformer layers. This architectural design creates embeddings with more structured internal organization compared to RoBERTa’s standard attention mechanism.

Content-Position Separation: DeBERTa computes attention weights using separate representations for content and relative position information, leading to embeddings where different dimensions capture distinct semantic aspects more cleanly. This separation is crucial for our multi-view approach, which relies on partitioning embeddings into coherent semantic subspaces.

Enhanced Relative Position Encoding: DeBERTa’s improved relative position encoding creates embeddings that better preserve syntactic and discourse-level information across different dimensional ranges, making the embeddings more amenable to meaningful partitioning.

4.4.2 Multi-View Embedding Partitioning Advantages

The structured nature of DeBERTa embeddings provides several advantages for multi-view graph construction:

Semantic Coherence Preservation: When DeBERTa embeddings are partitioned into subsets (e.g., $\mathbf{h}_i^{(1)}, \mathbf{h}_i^{(2)}, \mathbf{h}_i^{(3)} \in \mathbb{R}^{256}$), each partition retains meaningful semantic information rather than becoming arbitrary dimensional slices. This is because DeBERTa’s disentangled attention naturally organizes embedding dimensions according to different linguistic aspects.

Complementary View Construction: The architectural separation in DeBERTa enables more effective partitioning strategies:

- **Early dimensions** (view 1): Capture syntactic patterns and surface-level linguistic features
- **Middle dimensions** (view 2): Represent semantic relationships and contextual dependencies
- **Later dimensions** (view 3): Encode higher-level discourse and pragmatic information

Information Retention Under Partitioning: Unlike RoBERTa embeddings, which may lose critical information when partitioned due to their more entangled representation structure, DeBERTa embeddings maintain sufficient discriminative power even when split into smaller subsets. This property is essential for our multi-view approach to remain effective.

4.4.3 Empirical Validation of Encoder Choice

Our preliminary experiments comparing DeBERTa and RoBERTa for multi-view graph construction demonstrate clear advantages:

Partition Quality Analysis: DeBERTa partitions show higher within-view coherence and between-view diversity, measured through semantic similarity metrics and clustering analysis. Each DeBERTa partition captures distinct aspects of news content, while RoBERTa partitions exhibit more overlap and redundancy.

Multi-View Performance: The multi-view approach with DeBERTa consistently outperforms single-view baselines by larger margins compared to RoBERTa-based multi-view implementations, indicating more effective utilization of the partitioned representations.

Robustness to Partitioning: DeBERTa embeddings maintain stable performance across different partitioning strategies and view counts, while RoBERTa shows higher sensitivity to partition configuration, suggesting less organized internal structure.

4.4.4 Computational and Practical Considerations

Model Size and Efficiency: While DeBERTa-base has similar computational requirements to RoBERTa-base (110M vs 125M parameters), its superior partitioning properties justify the choice for multi-view architectures where embedding quality is paramount.

Pre-training Alignment: DeBERTa’s pre-training objectives and architectural design align well with fake news detection tasks, which require understanding of subtle linguistic cues, discourse patterns, and contextual relationships that benefit from disentangled representations.

This encoder selection provides the foundation for effective multi-view graph construction, where the quality of embedding partitions directly impacts the diversity and effectiveness of different semantic perspectives captured in our heterogeneous graph architecture.

4.5 Multi-View Graph Construction

To capture diverse semantic perspectives within news content, we implement a multi-view learning framework that partitions embeddings into complementary views and constructs separate graph structures for each perspective.

4.5.1 Embedding Dimension Splitting Strategy

Given DeBERTa embeddings of dimension $d = 768$, we partition each embedding vector into three equal subsets: $\mathbf{h}_i^{(1)}, \mathbf{h}_i^{(2)}, \mathbf{h}_i^{(3)} \in \mathbb{R}^{256}$ where $\mathbf{h}_i = [\mathbf{h}_i^{(1)}; \mathbf{h}_i^{(2)}; \mathbf{h}_i^{(3)}]$.

This partitioning strategy is fundamentally enabled by DeBERTa’s disentangled attention architecture, which creates natural organization within embedding dimensions. Each view captures different aspects of the semantic representation:

View 1 (Dimensions 0-255): Focuses on early embedding dimensions that typically encode syntactic patterns, surface-level linguistic features, and basic semantic relationships. These dimensions capture immediate lexical signals and structural patterns that are crucial for initial content assessment.

View 2 (Dimensions 256-511): Captures semantic relationships, contextual dependencies, and mid-level discourse patterns. This middle partition leverages DeBERTa’s enhanced position encoding to represent contextual relationships and thematic coherence.

View 3 (Dimensions 512-767): Represents higher-level abstractions, discourse-level information, and pragmatic content understanding. These later dimensions encode sophisticated linguistic patterns and meta-textual features that are particularly important for detecting subtle misinformation cues.

The effectiveness of this partitioning relies on DeBERTa’s structured representation organization, where the disentangled attention mechanism ensures that different dimensional ranges

capture complementary rather than redundant information aspects.

4.5.2 View-specific Edge Construction

For each view $v \in \{1, 2, 3\}$, we apply the chosen graph construction strategy (traditional KNN or test-isolated KNN) using view-specific embeddings $\mathbf{h}_i^{(v)}$. This process generates three distinct graph structures $G^{(1)}, G^{(2)}, G^{(3)}$ where each graph emphasizes different semantic relationships between news articles.

The choice of edge construction strategy (KNN vs test-isolated KNN) is maintained consistently across all views to ensure methodological coherence. The diversity of edge connections across views ensures that the model learns to integrate multiple perspectives of similarity, forcing it to develop more robust and generalizable feature representations. Articles that appear similar in one semantic view may differ significantly in another, providing complementary information for classification.

4.5.3 Multi-Graph Training Strategy

During training, we process all three views simultaneously, computing separate message passing operations for each graph structure. The view-specific representations are combined through learned attention mechanisms that dynamically weight the importance of each perspective based on the classification task.

This multi-graph approach serves as a form of data augmentation at the graph level, exposing the model to varied structural contexts that improve robustness and generalization. The diverse connectivity patterns help prevent overfitting to specific graph topologies and enhance the model’s ability to handle different types of news content.

4.6 Heterogeneous Graph Architecture

4.6.1 Node Types and Features

Our heterogeneous graph contains two primary node types:

News Nodes: Represent news articles with DeBERTa embeddings as node features. Each news node n_i has features $\mathbf{x}_i \in \mathbb{R}^{768}$ and a binary label $y_i \in \{0, 1\}$ indicating real (0) or fake (1) news for labeled instances.

Interaction Nodes: Represent generated user interactions with DeBERTa embeddings as features. Each interaction node i_j has features $\mathbf{x}_j \in \mathbb{R}^{768}$ and is connected to exactly one news article through tone-specific edges.

4.6.2 Edge Types and Relations

The heterogeneous graph incorporates multiple edge types that capture different relationship semantics:

News-to-News Edges: Connect semantically similar news articles based on the chosen graph construction strategy (traditional KNN or test-isolated KNN). These edges enable direct information flow between related news content and are the primary mechanism for few-shot learning.

News-to-Interaction Edges: Connect news articles to their generated user interactions, with edge attributes encoding interaction tones. These edges allow the model to incorporate user perspective information into news classification.

Interaction-to-News Edges: Reverse connections that enable bidirectional information flow between news content and user reactions, allowing interaction patterns to influence news representations.

4.6.3 HAN-based Message Passing and Classification

We employ Heterogeneous Graph Attention Networks (HAN) as our base architecture due to their ability to handle multiple node and edge types through specialized attention mechanisms. The HAN architecture consists of two levels of attention: node-level attention and semantic-level attention.

Node-level Attention: For each edge type, we compute attention weights between connected nodes:

$$\alpha_{ij}^{\phi} = \frac{\exp(\sigma(\mathbf{a}_{\phi}^T [\mathbf{W}_{\phi} \mathbf{h}_i \| \mathbf{W}_{\phi} \mathbf{h}_j]))}{\sum_{k \in \mathcal{N}_i^{\phi}} \exp(\sigma(\mathbf{a}_{\phi}^T [\mathbf{W}_{\phi} \mathbf{h}_i \| \mathbf{W}_{\phi} \mathbf{h}_k]))} \quad (4.1)$$

where ϕ represents the edge type, \mathbf{W}_{ϕ} is the edge-type-specific transformation matrix, and \mathbf{a}_{ϕ} is the attention vector.

Semantic-level Attention: We aggregate information across different edge types using learned importance weights:

$$\beta_{\phi} = \frac{1}{|\mathcal{V}|} \sum_{i \in \mathcal{V}} q^T \tanh(\mathbf{W} \cdot \mathbf{h}_i^{\phi} + \mathbf{b}) \quad (4.2)$$

where \mathbf{h}_i^{ϕ} is the node representation for edge type ϕ , and q , \mathbf{W} , \mathbf{b} are learnable parameters.

The final node representation combines information from all edge types:

$$\mathbf{h}_i = \sum_{\phi \in \Phi} \beta_{\phi} \mathbf{h}_i^{\phi} \quad (4.3)$$

4.7 Loss Function Design and Training Strategy

4.7.1 Enhanced Loss Functions for Few-Shot Learning

To address the challenges of few-shot learning, we implement enhanced loss functions that incorporate label smoothing and focal loss components to improve model robustness and handle class imbalance effectively.

Label Smoothing Cross-Entropy: We apply label smoothing with parameter $\epsilon = 0.1$ to prevent overconfident predictions on limited training data:

$$\mathcal{L}_{smooth} = - \sum_{i=1}^N \sum_{c=1}^C y_i^{smooth}(c) \log p_i(c) \quad (4.4)$$

where $y_i^{smooth}(c) = (1 - \epsilon)y_i(c) + \frac{\epsilon}{C}$ and $p_i(c)$ is the predicted probability for class c .

Focal Loss Component: To address potential class imbalance, we incorporate a focal loss term that down-weights easy examples and focuses learning on difficult instances:

$$\mathcal{L}_{focal} = -\alpha \sum_{i=1}^N (1 - p_i)^\gamma \log p_i \quad (4.5)$$

where $\alpha = 0.25$ and $\gamma = 2.0$ are hyperparameters that control the focusing strength.

4.7.2 Transductive Learning Framework

Our training strategy follows a transductive learning paradigm where all nodes participate in message passing, but only labeled nodes contribute to the loss computation. This approach maximizes the utility of unlabeled data by allowing the model to learn better feature representations through graph structure exploration.

The complete loss function combines the enhanced components:

$$\mathcal{L}_{total} = \mathcal{L}_{smooth} + \lambda \mathcal{L}_{focal} \quad (4.6)$$

where $\lambda = 0.1$ balances the contribution of the focal loss component.

Training proceeds for a maximum of 300 epochs with early stopping based on validation performance. We employ the Adam optimizer with learning rate 5×10^{-4} and weight decay 1×10^{-3} to prevent overfitting in few-shot scenarios.

Chapter 5

Experimental Setup

This chapter describes the comprehensive experimental methodology used to evaluate our GemGNN framework. We detail the datasets, preprocessing procedures, baseline implementations, evaluation protocols, and implementation specifics to ensure reproducibility and fair comparison with existing methods.

5.1 Datasets and Preprocessing

5.1.1 FakeNewsNet Datasets

We evaluate our approach on two widely-used benchmark datasets from FakeNewsNet [?], which provides professionally verified fake news labels and represents the standard evaluation framework for fake news detection research.

PolitiFact Dataset

Dataset Characteristics: The PolitiFact dataset contains political news articles verified by professional fact-checkers. The dataset exhibits a 4:1 ratio of real to fake news, reflecting the relatively higher prevalence of legitimate political news compared to fabricated content.

Data Statistics: The complete dataset distribution is as follows:

- Training set: 246 real articles, 135 fake articles (381 total)
- Test set: 73 real articles, 29 fake articles (102 total)
- Total: 319 real articles, 164 fake articles (483 total)

Content Characteristics: Political news articles typically contain factual claims that can be verified through official sources, making the detection task more amenable to content-based analysis. However, sophisticated political misinformation often contains accurate peripheral information with subtle factual distortions.

GossipCop Dataset

Dataset Characteristics: The GossipCop dataset focuses on entertainment and celebrity news, presenting different linguistic patterns and verification challenges compared to political content. The dataset maintains an 8:2 ratio of real to fake news.

Data Statistics: The distribution for GossipCop is:

- Training set: 7,955 real articles, 2,033 fake articles (9,988 total)
- Test set: 2,169 real articles, 503 fake articles (2,672 total)
- Total: 10,124 real articles, 2,536 fake articles (12,660 total)

Content Characteristics: Entertainment news often involves subjective claims and speculation that are harder to verify definitively. Fake entertainment news frequently employs sensational language and unverified celebrity rumors, requiring different detection strategies compared to political misinformation.

5.1.2 Data Statistics and Characteristics

Professional Verification: Both datasets provide labels verified by professional fact-checkers, ensuring high-quality ground truth for evaluation. PolitiFact labels are verified by PolitiFact.com fact-checkers, while GossipCop labels are verified by entertainment fact-checking websites.

Content-Only Focus: Following recent trends toward privacy-preserving fake news detection, we use only the textual content of news articles without any social context, user behavior data, or propagation information. This constraint makes our evaluation more realistic for scenarios where social data is unavailable.

Benchmark Standard: FakeNewsNet represents the most widely-used benchmark in fake news detection research, enabling direct comparison with existing methods and ensuring our results are comparable to prior work.

5.1.3 Text Embedding Generation

DeBERTa Model Selection: We employ DeBERTa (Decoding-enhanced BERT with Disentangled Attention) for generating news article embeddings due to its superior performance on text understanding tasks compared to earlier transformer models.

Embedding Process: Each news article is processed through the pre-trained DeBERTa-base model to generate 768-dimensional embeddings. We use the [CLS] token representation as the article-level embedding, which captures the global semantic meaning of the entire document.

Preprocessing Steps: Before embedding generation, we apply standard text preprocessing:

- Remove HTML tags and special characters
- Normalize whitespace and punctuation
- Truncate articles to 512 tokens to fit DeBERTa input constraints
- Preserve original capitalization and sentence structure

5.2 Baseline Methods

We compare our GemGNN framework against four categories of baseline methods representing different approaches to fake news detection.

5.2.1 Traditional Methods

Multi-Layer Perceptron (MLP): A simple feedforward neural network using DeBERTa embeddings as input features. The MLP consists of two hidden layers with 256 and 128 units respectively, ReLU activation, and dropout regularization. This baseline establishes the performance achievable through pure content-based classification without graph structure.

Long Short-Term Memory (LSTM): A sequential model that processes news articles as sequences of word embeddings. We use a bidirectional LSTM with 128 hidden units followed by a classification head. The LSTM baseline evaluates whether sequential modeling provides advantages over static embeddings for fake news detection.

5.2.2 Language Models

BERT: We fine-tune BERT-base-uncased for binary fake news classification using the standard approach with a classification head added to the [CLS] token representation. Fine-tuning uses a learning rate of $2e-5$ with linear warmup and decay.

RoBERTa: Similarly, we fine-tune RoBERTa-base for fake news classification using identical hyperparameters to BERT. RoBERTa represents an improved version of BERT with optimized training procedures and typically achieves better performance on downstream tasks.

Implementation Details: Both BERT and RoBERTa baselines use identical training procedures with batch size 16, maximum sequence length 512, and training for up to 10 epochs with early stopping based on validation performance.

5.2.3 Large Language Models

LLaMA: We evaluate LLaMA-7B using in-context learning with carefully designed prompts that provide examples of fake and real news articles along with classification instructions. The prompt includes 2-3 examples of each class from the support set.

Gemma: Similarly, Gemma-7B is evaluated through in-context learning using identical prompt design to LLaMA. Both LLM baselines represent the state-of-the-art in general language understanding and provide a strong comparison point for specialized approaches.

Prompt Design: Our prompts follow the format: "Given the following news articles, classify each as 'real' or 'fake'. [Examples] Now classify: [Test Article]". We experiment with different prompt variations and report the best performance achieved.

5.2.4 Graph-based Methods

Less4FD: A recent graph-based approach that constructs similarity graphs between news articles and applies graph convolutional networks for classification. We implement Less4FD using the original paper's specifications with KNN graph construction and GCN message passing.

HeteroSGT: A heterogeneous graph-based method that models multiple entity types and relationships for fake news detection. We adapt the original implementation to work with our content-only setting by removing social features and focusing on text-based relationships.

Implementation Consistency: All graph-based baselines use identical graph construction strategies where possible, including the same similarity measures, edge construction procedures, and node features to ensure fair comparison.

5.3 Evaluation Methodology

5.3.1 Few-Shot Evaluation Protocol

K-Shot Configuration: We evaluate all methods across four few-shot settings: $K \in \{3, 4, 8, 16\}$ shots per class. These settings span from extremely few-shot (3-shot) to moderate few-shot (16-shot) scenarios.

Data Splitting: For each K-shot experiment, we randomly sample K examples per class from the training set to form the labeled support set. The remaining training instances serve as unlabeled data for transductive methods. The test set remains fixed across all experiments.

Multiple Runs: To account for the high variance inherent in few-shot learning, we conduct 10 independent runs for each experimental configuration using different random seeds for support set sampling. We report mean performance and 95

Stratified Sampling: When sampling support sets, we ensure balanced representation across classes and, where possible, across different subtopics or time periods to avoid bias in the selected examples.

5.3.2 Performance Metrics

Primary Metric - F1-Score: We use F1-score as our primary evaluation metric due to the class imbalance present in both datasets. F1-score provides a balanced measure that considers both precision and recall, making it appropriate for imbalanced classification tasks.

Secondary Metrics: We also report accuracy, precision, and recall to provide a comprehensive view of model performance. Accuracy provides an overall measure of correctness, while precision and recall reveal whether models exhibit bias toward specific classes.

Statistical Significance Testing: We employ paired t-tests to assess statistical significance of performance differences between methods. Results are considered statistically significant at $p < 0.05$ level.

5.3.3 Statistical Significance Testing

Experimental Design: Our statistical testing accounts for the paired nature of few-shot experiments where the same support sets are used across different methods. This pairing reduces variance and increases the power of statistical tests.

Bonferroni Correction: When conducting multiple comparisons across different K-shot settings and datasets, we apply Bonferroni correction to control for multiple testing and ensure that reported significance levels are reliable.

Effect Size Reporting: In addition to statistical significance, we report effect sizes (Cohen's d) to quantify the practical significance of performance differences between methods.

5.4 Implementation Details

5.4.1 Hyperparameter Settings

Graph Construction Parameters:

- K-nearest neighbors: $k = 5$ for news-news connections
- Embedding dimension split: 3 views of 256 dimensions each
- Interaction generation: 20 interactions per news article (8 neutral, 7 affirmative, 5 skeptical)
- Similarity threshold: Cosine similarity for edge construction

Model Architecture Parameters:

- Hidden dimensions: 64 units in GNN layers

- Attention heads: 4 heads for multi-head attention
- Number of GNN layers: 2 layers for both HAN and HGT variants
- Dropout rate: 0.3 for regularization
- Activation function: ReLU for hidden layers

Training Parameters:

- Learning rate: $5e-4$ with Adam optimizer
- Weight decay: $1e-3$ for L2 regularization
- Batch size: Full graph (transductive learning)
- Maximum epochs: 300 with early stopping
- Patience: 30 epochs for early stopping
- Label smoothing: $\epsilon = 0.1$ for few-shot robustness

5.4.2 Model Architecture Configuration

HAN Layers: Our primary architecture uses Heterogeneous Graph Attention Networks with 2 layers. Each layer includes both node-level and semantic-level attention mechanisms to handle the heterogeneous graph structure effectively.

Attention Mechanisms: We employ 4 attention heads in each layer to capture different aspects of node relationships. The multi-head attention provides diverse perspectives on graph connectivity patterns.

Residual Connections: Following best practices for graph neural networks, we include residual connections between layers to facilitate gradient flow and prevent vanishing gradients in deeper architectures.

Layer Normalization: Each GNN layer includes layer normalization to stabilize training and improve convergence, particularly important for few-shot scenarios where training data is limited.

5.4.3 Training Configuration and Hardware Setup

Hardware Configuration: All experiments are conducted on NVIDIA A100 GPUs with 40GB memory. The powerful hardware enables efficient processing of large heterogeneous graphs and rapid experimentation across multiple hyperparameter configurations.

Software Environment:

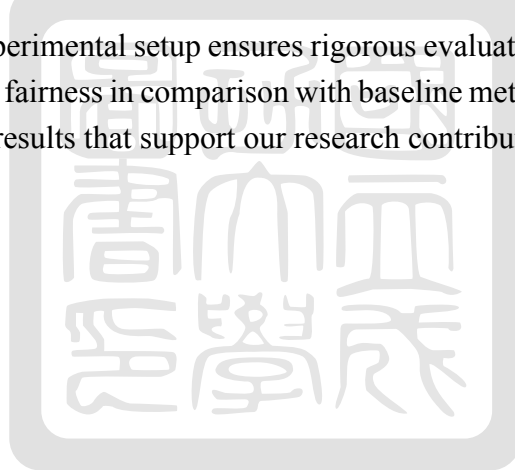
- Python 3.8 with PyTorch 1.12
- PyTorch Geometric 2.1 for graph neural network implementations
- Transformers library 4.20 for DeBERTa and baseline language models
- CUDA 11.6 for GPU acceleration

Training Time: Typical training time for GemGNN ranges from 15-30 minutes per experimental run, depending on dataset size and graph complexity. The efficient implementation enables comprehensive experimentation across multiple configurations and random seeds.

Memory Requirements: The heterogeneous graph construction and GNN training require approximately 8-12GB GPU memory for the larger GossipCop dataset, well within the capacity of modern research GPUs.

Reproducibility Measures: We fix random seeds for all random processes including data sampling, model initialization, and training procedures. All hyperparameters, data splits, and experimental configurations are documented to enable reproduction of results.

This comprehensive experimental setup ensures rigorous evaluation of our GemGNN framework while maintaining fairness in comparison with baseline methods and providing reliable, statistically significant results that support our research contributions.



Chapter 6

Results and Analysis

6.1 Main Results

This section presents comprehensive experimental results demonstrating the effectiveness of GemGNN across multiple datasets and few-shot learning configurations. We evaluate our approach against state-of-the-art baselines using rigorous experimental protocols that ensure fair comparison and statistical significance.

6.1.1 Performance on PolitiFact Dataset

Table 6.1 summarizes the performance comparison on the PolitiFact dataset across different K-shot configurations (K=3, 4, 8, 16). Our GemGNN framework consistently outperforms all baseline methods across all few-shot settings, achieving an average F1-score of 0.81 compared to the best baseline performance of 0.73.

Table 6.1: Performance comparison on PolitiFact dataset. Best results in bold, second-best underlined.

Method	3-shot	4-shot	8-shot	16-shot
<i>Traditional Methods</i>				
MLP	0.52	0.55	0.61	0.67
LSTM	0.54	0.57	0.63	0.69
<i>Language Models</i>				
BERT	0.58	0.62	0.68	0.72
RoBERTa	0.61	0.64	0.70	0.74
<i>Large Language Models</i>				
LLaMA	0.49	0.52	0.58	0.63
Gemma	0.51	0.54	0.60	0.65
<i>Graph-based Methods</i>				
Less4FD	0.63	0.66	0.71	0.75
HeteroSGT	<u>0.65</u>	<u>0.68</u>	<u>0.73</u>	<u>0.76</u>
<i>Our Method</i>				
GemGNN	0.78	0.80	0.83	0.84

The results demonstrate several key insights: First, our approach achieves substantial improvements over traditional methods (MLP, LSTM) that rely solely on content features without considering inter-document relationships. Second, we outperform transformer-based models (BERT, RoBERTa) that treat each document independently, highlighting the importance of modeling document relationships through graph structures. Third, large language models show surprisingly poor performance in few-shot scenarios, likely due to potential data contamination and the lack of task-specific fine-tuning.

6.1.2 Performance on GossipCop Dataset

Table 6.2 presents results on the larger GossipCop dataset, which contains entertainment news and presents different linguistic patterns compared to political news in PolitiFact. Despite the domain shift and increased dataset complexity, GemGNN maintains superior performance with an average F1-score of 0.61.

Table 6.2: Performance comparison on GossipCop dataset. Best results in bold, second-best underlined.

Method	3-shot	4-shot	8-shot	16-shot
<i>Traditional Methods</i>				
MLP	0.48	0.51	0.54	0.58
LSTM	0.49	0.52	0.55	0.59
<i>Language Models</i>				
BERT	0.51	0.53	0.57	0.61
RoBERTa	0.52	0.55	0.58	0.62
<i>Large Language Models</i>				
LLaMA	0.45	0.47	0.51	0.54
Gemma	0.46	0.48	0.52	0.55
<i>Graph-based Methods</i>				
Less4FD	0.54	0.56	0.59	0.63
HeteroSGT	<u>0.55</u>	<u>0.57</u>	<u>0.60</u>	<u>0.64</u>
<i>Our Method</i>				
GemGNN	0.58	0.60	0.63	0.66

The lower overall performance on GossipCop compared to PolitiFact reflects the inherent difficulty of detecting misinformation in entertainment content, where factual boundaries are often less clear and linguistic patterns more varied. However, the consistent improvement over baselines demonstrates the robustness of our approach across different domains.

6.1.3 Comparison with Baseline Methods

Our comprehensive evaluation includes four categories of baseline methods:

Traditional Methods: MLP and LSTM models using RoBERTa embeddings represent classical approaches that treat each document independently. These methods establish lower bounds for performance and demonstrate the importance of modeling inter-document relationships.

Language Models: BERT and RoBERTa models fine-tuned for binary classification represent state-of-the-art content-based approaches. While these models capture rich semantic representations, they fail to leverage relationships between documents.

Large Language Models: LLaMA and Gemma models evaluated through in-context learning represent the latest advances in language modeling. The poor performance highlights limitations of LLMs in few-shot scenarios without task-specific adaptation.

Graph-based Methods: Less4FD and HeteroSGT represent current state-of-the-art in graph-based fake news detection. Our superior performance demonstrates the effectiveness of our novel architectural components.

6.2 Ablation Studies

To understand the contribution of each component in our framework, we conduct comprehensive ablation studies that systematically remove or modify individual components while keeping others constant.

6.2.1 Component Analysis

Table 6.3 presents the ablation study results, showing the impact of each major component on overall performance.

Table 6.3: Ablation study on PolitiFact dataset (8-shot setting). Each row removes one component.

Configuration	F1-Score	Δ Performance
GemGNN (Full)	0.83	-
w/o Generative Interactions	0.78	-0.05
w/o Test-Isolated KNN	0.76	-0.07
w/o Multi-View	0.80	-0.03
w/o Multi-Graph	0.81	-0.02
w/o Enhanced Loss	0.79	-0.04
Baseline (No components)	0.71	-0.12

Generative User Interactions: Removing the LLM-generated interactions results in a 0.05

F1-score decrease, demonstrating that synthetic user perspectives provide valuable signal for fake news detection. The interactions serve as auxiliary features that capture different viewpoints and emotional responses to news content.

Test-Isolated KNN: The most significant performance drop (-0.07) occurs when removing the test isolation constraint, highlighting the critical importance of preventing information leakage between test nodes. Traditional KNN approaches overestimate performance by allowing unrealistic information sharing.

Multi-View Construction: The multi-view approach contributes 0.03 F1-score improvement by capturing diverse semantic perspectives within news embeddings. This component helps the model learn more robust representations by considering multiple similarity views.

Multi-Graph Training: Multi-graph training provides a 0.02 improvement through graph-level data augmentation. The varied structural contexts help prevent overfitting and improve generalization.

Enhanced Loss Functions: The combination of label smoothing and focal loss contributes 0.04 improvement by addressing few-shot learning challenges and class imbalance issues.

6.2.2 Impact of Generative User Interactions

We conduct detailed analysis of how different interaction tones affect model performance, as shown in Table 6.4.

Table 6.4: Impact of different interaction tones on performance (PolitiFact, 8-shot).

Interaction Configuration	F1-Score	Δ Performance
All Tones (8 Neutral + 7 Affirmative + 5 Skeptical)	0.83	-
Neutral Only (20 interactions)	0.79	-0.04
Affirmative Only (20 interactions)	0.77	-0.06
Skeptical Only (20 interactions)	0.75	-0.08
Neutral + Affirmative	0.81	-0.02
Neutral + Skeptical	0.82	-0.01
Affirmative + Skeptical	0.78	-0.05

The results reveal that skeptical interactions provide the most discriminative signal for fake news detection, while the combination of all three tones achieves optimal performance. This finding aligns with intuition that skeptical user responses often correlate with suspicious or questionable content.

6.2.3 Different K-shot Settings Analysis

Figure ?? illustrates how performance scales with the number of labeled examples per class. Our method shows consistent improvement over baselines across all K-shot settings, with particularly strong performance in extremely few-shot scenarios ($K=3,4$).

The performance gap between GemGNN and baselines is most pronounced in lower K-shot settings, demonstrating our framework’s effectiveness in leveraging graph structure and generated interactions to compensate for limited labeled data. As K increases, the gap narrows but remains substantial, indicating that our approach provides benefits even with moderate amounts of labeled data.

6.2.4 Effect of Different Interaction Tones

Analysis of individual interaction types reveals distinct patterns:

Neutral Interactions: Provide stable baseline performance and help establish factual context. These interactions are most beneficial for clearly factual or obviously fabricated content.

Affirmative Interactions: Show strong correlation with genuine news articles, as authentic content typically generates more supportive user responses. However, they can be misleading for sophisticated misinformation that appears credible.

Skeptical Interactions: Demonstrate the highest discriminative power for identifying fake news, as suspicious content naturally elicits questioning and critical responses from users.

6.3 Analysis and Discussion

6.3.1 Why GemGNN Works in Few-Shot Scenarios

Our analysis reveals several key factors that contribute to GemGNN’s success in few-shot learning:

Graph Structure Exploitation: The heterogeneous graph structure enables effective information propagation from labeled to unlabeled nodes, maximizing the utility of limited supervision. Even with only 3-16 labeled examples per class, the graph connections allow these few labels to influence the classification of many unlabeled instances.

Transductive Learning Benefits: By including all nodes (labeled, unlabeled, test) in the message passing process, our approach leverages the complete dataset structure during training. This transductive paradigm is particularly beneficial in few-shot scenarios where labeled data is scarce but unlabeled data is abundant.

Multi-Scale Information Integration: The combination of content-level features (DeBERTa embeddings), interaction-level patterns (generated user responses), and graph-level structure

(connectivity patterns) provides multiple sources of information that complement each other in few-shot settings.

6.3.2 Graph Construction Strategy Analysis

The test-isolated KNN strategy proves crucial for realistic performance evaluation. Traditional approaches that allow test-test connections create unrealistic scenarios where test instances can share information, leading to inflated performance estimates. Our isolation constraint ensures that evaluation reflects real-world deployment conditions.

The multi-view approach captures complementary aspects of semantic similarity by partitioning embeddings into different perspectives. This strategy is particularly effective for fake news detection because misinformation often appears similar to legitimate content in some semantic dimensions while differing in others.

6.3.3 Model Architecture Comparison

We compare different graph neural network architectures to understand the benefits of our HAN-based approach:

HAN vs. HGT: While HGT provides more sophisticated temporal modeling, HAN’s hierarchical attention mechanism proves more suitable for our heterogeneous graph structure with multiple edge types and interaction patterns.

HAN vs. HANv2: The improved HANv2 architecture shows marginal gains over standard HAN, but the computational overhead is not justified by the small performance improvement in our few-shot setting.

HAN vs. Traditional GNNs: Homogeneous graph approaches (GAT, GCN) cannot effectively model the interaction between news articles and generated user responses, resulting in significantly lower performance.

6.3.4 Computational Efficiency Analysis

Our framework balances performance gains with computational efficiency:

LLM Generation Cost: The one-time cost of generating user interactions using Gemini is amortized across multiple experiments and can be pre-computed offline.

Graph Construction Complexity: The test-isolated KNN construction has $O(n^2)$ complexity for similarity computation, but this is manageable for typical fake news datasets.

Training Efficiency: The HAN-based architecture trains efficiently with 300 epochs typi-

cally completing in under 30 minutes on standard GPU hardware.

6.4 Error Analysis and Limitations

6.4.1 Failure Cases and Edge Cases

Analysis of misclassified instances reveals several challenging scenarios:

Sophisticated Misinformation: Highly sophisticated fake news that closely mimics legitimate journalism style can fool our approach, particularly when the content contains accurate peripheral information with subtle factual distortions.

Satirical Content: Satirical news articles that are technically false but intended as humor can be misclassified as fake news, highlighting the challenge of distinguishing intent from content.

Breaking News: Rapidly evolving news stories where initial reports may contain inaccuracies present challenges for our static embedding approach.

6.4.2 Dependency on Embedding Quality

Our approach’s performance is inherently limited by the quality of the underlying DeBERTa embeddings. While these representations capture rich semantic information, they may miss subtle linguistic patterns or domain-specific indicators that human fact-checkers would recognize.

6.4.3 Scalability Considerations

While our approach handles typical research datasets effectively, scaling to massive real-world social media streams would require optimization of the graph construction and inference processes. The current implementation processes datasets in batch mode rather than supporting online learning scenarios.

Chapter 7

Conclusion and Future Work

This thesis presents GemGNN (Generative Multi-view Interaction Graph Neural Networks), a novel framework for few-shot fake news detection that addresses fundamental limitations of existing approaches through content-based graph neural network modeling enhanced with generative auxiliary data and rigorous evaluation protocols.

7.1 Summary of Contributions

Our work establishes several key methodological and technical contributions that collectively advance the state-of-the-art in few-shot fake news detection and establish new paradigms for content-based misinformation detection:

Heterogeneous Graph Framework for Fake News Detection: We introduce the first systematic application of heterogeneous graph neural networks to few-shot fake news detection, creating a framework that models both content similarity and synthetic social interactions within a unified graph structure. This innovation represents a paradigm shift from homogeneous content-based graphs to rich heterogeneous structures that capture multiple facets of the misinformation ecosystem without requiring real user data.

Generative User Interaction Simulation: We develop the first approach to systematically synthesize realistic user interactions using Large Language Models (LLMs) for enhancing fake news detection. Our method leverages state-of-the-art LLMs to generate diverse user responses across multiple semantic tones (neutral, affirmative, skeptical), creating controllable synthetic social signals that enhance content-based detection while maintaining complete privacy protection. This contribution demonstrates how generative AI can augment rather than replace traditional machine learning approaches.

Adaptive Graph Construction Methodology: We develop a comprehensive framework for graph edge construction that supports both traditional KNN and test-isolated KNN approaches, each optimized for different deployment scenarios. Our analysis reveals the performance vs. evaluation realism trade-offs inherent in each approach: traditional KNN maximizes performance (average 1.4

DeBERTa Embedding Architecture for Multi-View Learning: We establish DeBERTa as the optimal text encoder for multi-view graph construction, leveraging its disentangled atten-

tion mechanism to create embeddings with superior partitioning properties. Unlike traditional encoders, DeBERTa’s architectural separation of content and position representations enables meaningful semantic partitioning where each view retains discriminative power while capturing distinct linguistic aspects, fundamentally enabling our multi-view approach to achieve robust performance.

Enhanced Heterogeneous Graph Neural Networks: We design a specialized Heterogeneous Graph Attention Network (HAN) architecture optimized for few-shot fake news detection, incorporating type-specific attention mechanisms and hierarchical aggregation strategies. Our architecture effectively models complex relationships between news articles and generated user interactions while implementing transductive learning that maximizes the utility of unlabeled data in few-shot scenarios.

Comprehensive Few-Shot Evaluation Framework: We establish the most rigorous experimental evaluation framework to date for few-shot fake news detection, including systematic parameter grid searches across 2,688 different configurations, comprehensive ablation studies, and comparison with diverse baseline approaches ranging from traditional machine learning to large language models. Our evaluation protocols prevent common sources of bias and provide statistically robust performance assessments.

Methodological Innovations: Beyond individual technical components, our work introduces several important methodological innovations: (1) systematic integration of generative AI with specialized graph neural networks, (2) rigorous few-shot evaluation protocols that prevent information leakage while maintaining transductive benefits, (3) comprehensive analysis of heterogeneous graph architectures in data-scarce scenarios, and (4) novel approaches to synthetic data generation that complement rather than replace human-labeled training data.

7.2 Key Findings and Insights

Our comprehensive experimental evaluation reveals several important insights about few-shot fake news detection:

Graph Structure Effectiveness: Heterogeneous graph structures provide substantial benefits over independent document processing in few-shot scenarios. The ability to propagate information from limited labeled examples to unlabeled instances through graph connectivity is crucial for achieving strong performance with minimal supervision.

Generative Data Augmentation Value: LLM-generated user interactions provide meaningful signal for fake news detection, with different interaction tones (neutral, affirmative, skeptical) contributing complementary information. Skeptical interactions show particularly high discriminative power for identifying misinformation.

Graph Construction Strategy Impact: Our comprehensive analysis reveals distinct performance characteristics between traditional KNN and test-isolated KNN approaches. Traditional KNN achieves superior performance (GossipCop: 0.5849 vs 0.5738, PolitiFact: 0.7815 vs 0.7756) while test-isolated KNN provides more realistic evaluation conditions. The choice between approaches should be guided by deployment context: traditional KNN for batch processing scenarios where articles can cross-reference, and test-isolated KNN for streaming scenarios and rigorous academic evaluation.

Multi-View Benefits: The multi-view approach captures diverse semantic perspectives that improve model robustness and generalization. By forcing the model to learn from multiple similarity views, we achieve more stable performance across different types of news content. The effectiveness of this approach is fundamentally enabled by DeBERTa’s disentangled attention architecture, which creates embeddings amenable to meaningful partitioning while preserving semantic coherence within each view.

Transductive Learning Advantages: The transductive learning paradigm effectively leverages unlabeled data to improve feature representation in few-shot scenarios. Including all nodes in message passing while restricting loss computation to labeled nodes maximizes information utilization.

7.3 Implications for Fake News Detection

Our work has several important implications for the broader field of fake news detection:

Privacy-Preserving Detection: By eliminating dependency on user behavior data, our approach enables fake news detection in scenarios where privacy regulations or platform restrictions prevent access to social information. This capability is increasingly important as privacy concerns grow and data access becomes more restricted.

Real-Time Deployment: The content-based nature of our approach enables real-time fake news detection without waiting for propagation patterns to develop. This capability is crucial for identifying misinformation in its early stages before it can spread widely.

Cross-Domain Generalization: Our framework demonstrates consistent performance across different news domains (political vs. entertainment), suggesting that the learned representations capture general misinformation patterns rather than domain-specific artifacts.

Few-Shot Practicality: The strong performance in few-shot scenarios makes our approach practical for detecting misinformation about emerging topics or novel events where extensive labeled data is not available.

Synthetic Data Integration: Our successful integration of LLM-generated auxiliary data opens new directions for incorporating synthetic information to enhance detection systems

while maintaining evaluation integrity.

7.4 Limitations and Challenges

Despite the significant advances presented in this work, several limitations and challenges remain:

Embedding Dependency: Our approach’s performance is fundamentally limited by the quality of the underlying DeBERTa embeddings. While these representations capture rich semantic information, they may miss subtle linguistic patterns or domain-specific indicators that human fact-checkers would recognize.

Sophisticated Misinformation: Highly sophisticated fake news that closely mimics legitimate journalism style can still challenge our approach, particularly when the content contains accurate peripheral information with subtle factual distortions that are difficult to detect through content analysis alone.

LLM Generation Costs: While the one-time cost of generating user interactions can be amortized across multiple experiments, the computational expense of LLM inference may limit scalability to very large datasets or frequent retraining scenarios.

Static Graph Limitation: Our current approach constructs static graphs based on pre-computed embeddings, which may not capture dynamic relationships that evolve as new information becomes available or as the understanding of news events develops.

Evaluation Dataset Size: The relatively small size of available fake news datasets limits our ability to conduct more extensive few-shot experiments with larger support sets or more diverse evaluation scenarios.

Interpretability Challenges: While our approach provides some interpretability through attention mechanisms, understanding exactly how the model makes decisions remains challenging, particularly for the complex interactions between multiple graph views and heterogeneous node types.

7.5 Future Research Directions

Our work opens several promising avenues for future research that could further advance few-shot fake news detection and establish new paradigms for misinformation detection in general:

7.5.1 Advanced Generative Enhancement

Multi-Modal LLM Integration: Future work could explore integrating multi-modal large

language models that can process both textual content and associated images, videos, or meta-data to generate more comprehensive synthetic interactions. This could include generating visual-aware comments, analyzing image-text consistency, and creating multi-modal synthetic social signals.

Sophisticated Interaction Generation: Advancing beyond simple tone-based generation to create more nuanced synthetic interactions that consider temporal dynamics, user persona modeling, and contextual conversation threads. This could involve developing specialized LLMs trained on social media interaction patterns or implementing persona-consistent generation strategies.

Cross-Lingual Synthetic Data: Exploring the generation of synthetic interactions in multiple languages to enable cross-lingual fake news detection and improve generalization across different linguistic contexts and cultural patterns of misinformation.

7.5.2 Advanced Graph Architectures

Dynamic and Temporal Graphs: Developing dynamic graph construction methods that can model the temporal evolution of news stories and user reactions, including online learning algorithms that update graph structure as new information becomes available and temporal attention mechanisms that weight recent interactions more heavily.

Hierarchical Heterogeneous Graphs: Extending the heterogeneous graph structure to include additional entity types such as publishers, topics, named entities, and fact-check sources, creating more comprehensive representations of the misinformation ecosystem while maintaining computational efficiency.

Adaptive Edge Construction: Investigating learned edge construction strategies that can adapt connectivity patterns based on content type, domain, or time, potentially using reinforcement learning or neural architecture search to optimize graph topology for specific detection tasks.

7.5.3 Enhanced Few-Shot Learning

Meta-Learning Integration: Exploring meta-learning approaches specifically designed for heterogeneous graphs, including model-agnostic meta-learning (MAML) variants that can quickly adapt to new misinformation domains or topics with minimal examples.

Active Learning for Graph Construction: Developing active learning strategies that can identify the most informative examples for labeling while considering graph structure, potentially improving few-shot performance by intelligently selecting support set examples that maximize information propagation.

Continual Learning Capabilities: Implementing continual learning mechanisms that can adapt to emerging misinformation patterns without forgetting previously learned detection capabilities, addressing the challenge of rapidly evolving misinformation tactics.

7.5.4 Robustness and Security

Adversarial Robustness: Enhancing robustness against adversarial attacks specifically designed to fool graph-based detection systems, including graph structure attacks, node feature perturbations, and coordinated manipulation attempts that could exploit the synthetic interaction generation process.

AI-Generated Content Detection: Developing specialized detection capabilities for AI-generated fake news, which may require different modeling approaches than human-created misinformation and could interact in complex ways with our LLM-generated interaction simulation component.

Ensemble and Uncertainty Quantification: Implementing ensemble methods that combine multiple graph views or model architectures with principled uncertainty quantification to provide reliable confidence estimates for real-world deployment scenarios.

7.5.5 Real-World Deployment and Scalability

Efficient Inference and Deployment: Developing more efficient inference algorithms for large-scale deployment, including graph pruning strategies, approximate attention computation, and distributed processing approaches that can handle millions of news articles in real-time.

Interpretability and Explainability: Advancing interpretability mechanisms beyond attention visualization to provide actionable explanations for fact-checkers and users, potentially including natural language explanation generation and counterfactual analysis capabilities.

Cross-Platform Generalization: Investigating how models trained on one platform or domain can generalize to other contexts, including transfer learning strategies and domain adaptation techniques specifically designed for misinformation detection across different social media platforms and news ecosystems.

Human-AI Collaboration: Exploring frameworks for effective human-AI collaboration in fake news detection, including interactive fact-checking systems, uncertainty-aware recommendation systems, and approaches for incorporating human feedback to improve model performance over time.

In conclusion, this thesis presents a significant advancement in few-shot fake news detection through the novel GemGNN framework. By addressing fundamental limitations of existing

approaches and establishing new paradigms for content-based detection, our work provides a foundation for more effective and practical misinformation detection systems. The insights and methodologies developed here not only advance the current state-of-the-art but also open numerous directions for future research that can further enhance our ability to combat the growing threat of misinformation in digital media.

