LLM-based Contrastive Self-Supervised AMR Learning with Masked Graph Autoencoders for Fake News Detection

Anonymous authors
Paper under double-blind review

Abstract

The proliferation of misinformation in the digital age has led to significant societal challenges. Existing approaches often struggle with capturing long-range dependencies, complex semantic relations, and the social dynamics influencing news dissemination. Furthermore, these methods require extensive labelled datasets, making their deployment resource-intensive. In this study, we propose a novel self-supervised misinformation detection framework that integrates both complex semantic relations using Abstract Meaning Representation (AMR) and news propagation dynamics. We introduce an LLM-based graph contrastive loss (LGCL) that utilizes negative anchor points generated by a Large Language Model (LLM) to enhance feature separability in a zero-shot manner. To incorporate social context, we employ a multi-view graph masked autoencoder, which learns news propagation features from social context graph. By combining these semantic and propagation-based features, our approach effectively differentiates between fake and real news in a self-supervised manner. Extensive experiments demonstrate that our self-supervised framework achieves superior performance compared to other state-of-the-art methodologies, even with limited labelled datasets while improving generalizability.¹

1 Introduction

The spread of misinformation has become a significant problem in the digital age. It can lead to social unrest, foster hatred, erode trust, and ultimately impede the overall progress and stability of the society (Dewatana & Adillah, 2021). Hence, effectively detecting misinformation has become an essential challenge to solve. Yin et al. (2008) introduced the concept of 'veracity problem on the web' by designing a solution called TruthFinder. The method verified news content by cross-referencing it with information from reputable websites. Later, Feng et al. (2012) employed manually crafted textual features for detecting misinformation. However, manually crafted features are time-consuming to create and fail to capture the complex semantic relations present in the text. Subsequently, many researchers turned to more advanced techniques, utilizing RNN's, and Transformer-based (Long et al., 2017; Liu & Wu, 2018) models to address this issue. For example, RNNs are employed to capture local and temporal dependencies within text data (Ma et al., 2016b; Li et al., 2022) and BERT has been increasingly utilized to improve the comprehension of contextual relationships in news articles Devlin et al. (2019). Key limitations of these approaches are their struggle to maintain longer text dependencies and they do not capture complex semantic relations, such as events, locations, and trigger words. Gupta et al. (2025) used semantic relations through Abstract Meaning Representation (AMR) graph to solve this problem but their method requires label data for supervision. Additionally, these models often neglect the social context and dynamics that influence news propagation Yuan et al. (2019). Acknowledging this, researchers have introduced graph-based approaches that integrate social context referred as Social Context Graph (SCG) into the detection process (Min et al., 2022; Sun et al., 2022; Li et al., 2024). Despite their effectiveness, these methods rely heavily on large, labelled datasets for training. Collecting and annotating such extensive datasets is time-consuming and resource-intensive, limiting their practical implementation. To address this Yin et al. (2024) propose a Graph Masked Autoencoder with augmentations (GMA²) based model to generate unsupervised features from the social context graph but

¹Code repository: https://anonymous.4open.science/r/Fake1-3245/README.md

Table 1: Comparison of different methods based on their utilization of various graph-based learning components.

Method	AMR	SCG	GMA^2	GMA ² +Remasking	Unsupervised
$\overline{\rm EA^2N}$	✓	Х	Х	Х	Х
GACL	X	✓	X	X	×
$(\mathrm{UMD})^2$	X	✓	X	X	✓
GTUT	X	✓	X	X	✓
GAMC	X	✓	✓	X	✓
Ours	1	✓	✓	✓	✓

do not consider the semantic relationship within the text. Therefore, we require a model that is capable of incorporating semantic text features, a social context propagation graph and also perform well with minimal labelled data as highlighted in Table 1.

This paper proposes a novel self-supervised misinformation detection methodology that considers complex semantic relations among entities in the news and the propagation of the news as a social context graph. In order to identify the semantic relations, this method incorporates a self-supervised AMR encoder using the proposed graph contrastive loss. This loss helps in increasing the separation between fake and real classes in the latent space by LLM guided negative anchor points. In order to integrate the social context and capture the propagation of the news, our methodology also integrates a multi-view Graph Masked Autoencoder that employs the context and content of the news propagation process as the self-supervised signal to enhance the final feature space. These features, even with limited labelled data, achieve performance comparable or better than supervised counterparts using a simple linear SVM layer. The key contributions of our research are as follows:

- A novel self-supervised learning based on AMR and social context graph is introduced in order to validate the veracity of news articles, eliminating dependence on labelled data.
- In order to segregate the feature space among real and fake classes, graph contrastive loss is proposed. An LLM-based negative sampler is designed to handle negatives in the loss.
- To capture the social context and propagation feature of the news, we propose an augmentation-based multi-view masked graph autoencoder with remasking module.
- Comprehensive evaluation with SOTA methods, demonstrating its superior performance.

2 Related Work

In this section, we provide a concise overview of the approaches utilized for detecting misinformation. The relevant studies are categorized into two main components: misinformation detection and self-supervised graph learning methodologies.

2.1 Misinformation Detection Methods

Early research on misinformation detection focused on manually crafted linguistic features (Feng et al., 2012; Ma et al., 2016a; Long et al., 2017), requiring significant effort for evaluation. EANN (Wang et al., 2018) is proposed to effectively extract event-invariant features from multimedia content, thereby enhancing the detection of misinformation on newly arrived events. In this line of work, FakeFlow (Ghanem et al., 2021) classified news using lexical features and affective information. In a separate line of work, external knowledge was integrated to improve model performance. Different source of external knowledge was used, for example, Popat et al. (2017) retrieved external articles to model interactions; KAN (Dun et al., 2021) and CompareNet (Hu et al., 2021) leveraged Wikidata for domain expansion, while KGML (Yao et al., 2021) bridged meta-training and meta-testing using knowledge bases. Further, researchers have developed graph-based methods that incorporate social context into the detection process, for example, authors of

GTUT (Gangireddy et al., 2020) construct a graph for initial fake news spreader identification, (UMD)² (Silva et al., 2024) considers user credibility and propagation speed, GACL (Sun et al., 2022) constructs a tree of tweets for contrastive learning. All these methods do not leverage the complete propagation graph, and GACL requires supervision. Other graph-based methods like Min et al. (2022); Li et al. (2024) rely heavily on manual annotation and external data.

Recently, Abstract Meaning Representation (AMR)-based methods emerged to mitigate long-text dependency. Abstract Meaning Representation (AMR), as introduced by Banarescu et al. (2013), captures relationships between nodes using PropBank framesets. Recently, Zhang et al. (2023) utilized AMR to detect out-of-context multimodal misinformation by identifying discrepancies between textual and visual data. In (Gupta et al., 2023), authors encoded textual information using AMR and explored how its semantic relations influence the veracity assessment of news. However, this study lacked sufficient evidence or justification for entity relationships within the AMR graph. Further, in the integration of evidence in AMR, EA²N (Gupta et al., 2025) is proposed that effectively captures evidence among entities present in AMR. All of these approaches rely on supervised data for and have not explored the potential of unsupervised methods.

2.2 Self-Supervised Graph Learning

Self-supervised graph learning harnesses the structural richness of graph data to derive meaningful representations without relying on explicit labels (Wu et al., 2023). A Graph Auto-Encoder (GAE) based model that learns low-dimensional graph representations is proposed in (Kipf & Welling, 2016). Later studies improved GAEs by focusing on reconstructing masked node features to enhance self-supervised learning for classification (Hou et al., 2022). Further, Hou et al. (2023) improved the performance by introducing multi-view random remasking. Recently, an unsupervised method for detecting misinformation GAMC (Yin et al., 2024) has been proposed by leveraging both the context and content of news propagation as self-supervised signals. However, GAMC does not effectively handle complex semantic relations for longer text dependencies.

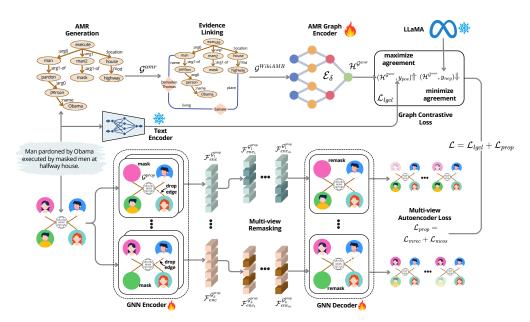


Figure 1: Overview of the proposed method: The news article is converted to an AMR graph \mathcal{G}^{amr} . \mathcal{G}^{amr} is then linked to external evidences from Wikipedia represented as $\mathcal{G}^{WikiAMR}$. This $\mathcal{G}^{WikiAMR}$ graph is then converted to latent space features $\mathcal{H}^{\mathcal{G}^{amr}}$ by the graph transformer \mathcal{E}_{δ} based on \mathcal{L}_{lgcl} optimization. The propagation graph of the same news article is then extracted and multiple augmentations are created. These augmented graphs are then passed to our multi-view remasked graph autoencoder which is optimized using \mathcal{L}_{prop} . The propagation graph feature $\mathcal{H}^{\mathcal{G}^{prop}}$ for each news is extracted from the trained GNN encoder. The final features for misinformation classification are obtained by concatenating $\mathcal{H}^{\mathcal{G}^{amr}}$ and $\mathcal{H}^{\mathcal{G}^{prop}}$.

3 Methodology

An overview of the proposed methodology is presented in Figure 1. In this section we present these in more detail.

3.1 Self-supervised AMR Graph Learning

Given an input text T, we first create the AMR graph $\mathcal{G}^{amr}(\mathcal{V}^{amr}, \mathcal{E}^{amr})$ capturing the relationships between different entities. AMR generation process involves parsing the sentences to extract linguistic information, including semantic roles, relations, and core events. In order to incorporate reasoning through AMR, we have integrated the external evidence by using the Evidence Linking Algorithm (ELA) used in Gupta et al. (2025). The graph after applying ELA is referred to as WikiAMR, represented as $\mathcal{G}^{WikiAMR}$. In the paper, authors have shown the importance of WikiAMR over AMR. WikiAMR comprises interconnected undirected paths between entity nodes in \mathcal{G}^{amr} generated from the text. The WikiAMR representation helps to distinguish the difference between real and fake articles.

AMR Graph Learning with Path Optimization: This module plays an important role in extracting meaningful features from the given WikiAMR graph. Features extracted here capture essential semantic relationships, enabling a deeper understanding of the underlying textual data. At the core of this module is a Graph Transformer (Cai & Lam, 2020), which employs various attention mechanisms to effectively process the graph representation. This allows the model to reason about and learn from the text more comprehensively.

The WikiAMR graph is first passed through a node initialization and relation encoder to transform it into a representation in $\mathbb{R}^{n \times k \times d}$, where n, k, and d denote the batch size, maximum sequence length, and the dimensionality of the graph encoding, respectively. To facilitate the model in identifying specific paths within $\mathcal{G}^{WikiAMR}$, the relation encoder computes the shortest path between two entities. This sequence of the path is subsequently converted into a relation vector using a Gated Recurrent Unit (GRU)-based RNN (Cho et al., 2014). q_t is the sequence encoding extracted from GRU to get the relation vector r_{uv} . The mathematical formulation for this encoding is given by:

$$\overrightarrow{q}_{t} = GRU_{f}(\overrightarrow{q}_{t-1}, sp_{t})$$

$$\overleftarrow{q}_{t} = GRU_{b}(\overleftarrow{q}_{t+1}, sp_{t})$$

Here, sp_t represents the shortest path between two entities. Formally, the shortest relation path $sp_{i\to j} = [e(u,k_1),e(k_1,k_2),\ldots,e(k_n,v)]$ between the node u and the node v, where $e(\cdot,\cdot)$ indicates the edge label and $k_{1:n}$ are the relay nodes. To compute the attention scores, the final relational encoding r_{uv} is split into two distinct components, $r_{u\to v}$ and $r_{v\to u}$, via a linear transformation with a parameter matrix W_r :

$$r_{uv} = [\overrightarrow{q}_n; \overleftarrow{q}_0], \quad [r_{u \to v}; r_{v \to u}] = W_r r_{uv}$$

Subsequently, attention scores β_{uv} are calculated by incorporating both entity and relation representations from the graph $\mathcal{G}^{WikiAMR}$:

$$\beta_{uv} = h(e_u, e_v, r_{uv})$$

$$= (e_u + r_{u \to v}) W_p^\top W_k (e_v + r_{v \to u})$$

$$= \underbrace{e_u W_p^\top W_k e_v}_{\textcircled{a}} + \underbrace{e_u W_p^\top W_k r_{v \to u}}_{\textcircled{b}}$$

$$+ \underbrace{r_{u \to v} W_p^\top W_k e_v}_{\textcircled{C}} + \underbrace{r_{u \to v} W_p^\top W_k r_{v \to u}}_{\textcircled{d}}$$

$$\tag{1}$$

The attention weights computed here guide the focus on entities according to their relationships. Each term in Equation 1 serves a distinct purpose: (a) models content-based attention, (b) captures biases related to the source of the relationship, (c) addresses biases from the target, and (d) encodes a general relational bias, providing a comprehensive view of entity interactions. Finally, the Graph Transformer (\mathcal{E}_{δ}) encodes $\mathcal{G}^{WikiAMR}$, producing the final graph representation as follows:

$$\mathcal{H}^{\mathcal{G}^{amr}} = \mathcal{E}_{\delta}(\mathcal{G}^{WikiAMR}) \in \mathbb{R}^{n \times k \times d}$$
(2)

Here, $\mathcal{H}^{\mathcal{G}^{amr}}$ represents the output graph embeddings generated by the Graph Transformer, and d is the feature dimensionality.

Graph Contrastive Loss: Our proposed LLM-guided graph contrastive loss (LGCL) function comprises two primary objectives. The first objective aims to ensure that the graph embedding remains close to its original embedding space by minimizing the reconstruction error between the predicted feature and the original feature. The second objective seeks to maximize the divergence between the predicted feature and the negative sample feature. To quantify the similarity between features, we utilize the Scaled Cosine Error (SCE) (Hou et al., 2022). Formally, given the original feature Y and the reconstructed output Y', SCE is defined as:

$$\mathcal{L}_{\text{SCE}} = \frac{1}{|\mathcal{N}|} \sum_{n \in \mathcal{N}} \left(1 - \frac{y_i^T y_i'}{\|y_i\| \cdot \|y_i'\|} \right)^{\gamma}, \quad \gamma \ge 1$$
 (3)

Here, γ is a scaling factor. When predictions have high confidence, the resulting cosine errors are generally less than 1 and diminish more quickly towards zero as the scaling factor $\gamma > 1$.

The contrastive loss requires both a positive sample feature y_{pos} and a negative sample feature y_{neg} to compare against the predicted feature. In the proposed formulation, $\mathcal{H}^{\mathcal{G}^{amr}}$ is used as y', y_{pos} is the original BERT-derived feature of the input text, while y_{neg} is a negative sample feature generated using an LLM-guided negative sampler. The final contrastive loss for graph-based self-supervised learning (SSL) is formulated as follows:

$$\mathcal{L}_{lgcl} = \mathcal{L}_{SCE}(y', y_{pos}) + \lambda \cdot \max(0, m - \mathcal{L}_{SCE}(y', y_{neg}))$$
(4)

Here, λ is a weighting factor, and m is the margin to ensure negatives are pushed apart in cosine space.

LLM-guided Negative Sampler: Let $\mathcal{X} = \{x_1, x_2, \dots, x_n\}$ denote the set of input features. For each input $x_i \in \mathcal{X}$, an LLM (zero-shot inference) assigns a pseudo label $\tilde{y}_i \in \{0, 1\}$, where:

$$\widetilde{y}_i = \begin{cases} 1 & \text{if } x_i \text{ is labelled as real,} \\ 0 & \text{if } x_i \text{ is labelled as fake.} \end{cases}$$

Using the LLM's output labels, we partition the input samples into two groups:

$$\mathcal{X}_{\text{real}} = \{x_i \mid \widetilde{y}_i = 1\}, \quad \mathcal{X}_{\text{fake}} = \{x_i \mid \widetilde{y}_i = 0\}.$$

We compute the centroids of the real and fake samples as,

$$c_{\mathrm{real}} = \frac{1}{|\mathcal{X}_{\mathrm{real}}|} \sum_{x_i \in \mathcal{X}_{\mathrm{real}}} \mathbf{f}_i, \quad c_{\mathrm{fake}} = \frac{1}{|\mathcal{X}_{\mathrm{fake}}|} \sum_{x_i \in \mathcal{X}_{\mathrm{fake}}} \mathbf{f}_i.$$

where a feature vector $\mathbf{f}_i \in \mathbb{R}^{n \times k \times d}$ is the initial BERT feature corresponding to x_i . We use c_{fake} as the representative negative sample for the real input sample, while c_{real} is used as the negative sample for the

fake input sample. The negative samples (y_{neg}) thus chosen is used in the contrastive loss of Equation 4. Specifically, LLaMA3-7B is used to generate these negative samples (y_{neg}) . Note that the pseudo labels generated by the LLM are noisy. However, the centroid thus identified shows stability in the practical experiments. Further note that no other pseudo labels except the centroid is used in the main pipeline.

LLM's Zero Shot Input Prompt:

Write in one word among 'real' or 'fake' whether given text is real or fake. {text}

LLM's Output: fake/real

3.2 Multi-View Social Context and Propagation Graph Learning

Each news article is converted into a propagation graph $G^{prop} = (V, E, \mathcal{F})$ as in Dou et al. (2021). Nodes in V represent one news article and users who forward that article. An edge in E exists between two nodes if there exists a forwarding relationship between them. The features for the news node are generated by passing the news article to a pre-trained language model (BERT), and the features for the user nodes are generated based on their recent 200 posts. The news and user node features are collectively referred to as \mathcal{F} .

Graph Augmentation: We use two augmentation strategies: ① feature masking and ② random edge removal for creating augmentations of the input graph as suggested in Yin et al. (2024). For input feature masking, we randomly select 50% nodes in the graph and replace their features with a masked token. For ②, we randomly remove 20% edges from the graph. Each augmented graph for G^{prop} is denoted as G^{prop} .

Graph Encoding: We encode each \mathcal{G}_i^{prop} into a latent space representation using a GNN encoder. For this, we use GIN (Xu et al., 2019) represented using Equation 5 as it is theoretically proven to distinguish between graph structures.

$$f_v^{(k)} = \text{MLP}\left((1+\epsilon) \cdot f_v^{(k-1)} + \sum_{u \in \mathcal{N}(v)} f_u^{(k-1)}\right)$$
 (5)

Here, $f_v^{(k)}$ is embedding of node v at layer k, $\mathcal{N}(v)$ contains neighbors of node v and ϵ is a learnable scalar controlling residual connections. The final node embeddings from the encoder for each \mathcal{G}_i^{prop} is represented as $\mathcal{F}_{enc}^{\mathcal{G}_i^{prop}}$.

For downstream classification tasks on G^{prop} we use the graph embedding $\mathcal{H}^{G^{prop}}$ calculated as:

$$\mathcal{H}^{G^{prop}} = \frac{1}{|V|} \sum_{v \in V} f_v \in \mathcal{F}_{enc}^{G^{prop}} \tag{6}$$

Multi-View Graph Decoding: Now, from the encoded node representations $\mathcal{F}_{enc}^{\mathcal{G}_{i}^{prop}}$, we decode the input node features \mathcal{F} using GIN as a decoder. In Yin et al. (2024) the authors use a single stage remasking for each $\mathcal{F}_{enc}^{\mathcal{G}_{i}^{prop}}$ to reconstruct the input features. But authors in Hou et al. (2023) have shown that feature reconstruction is susceptible to congruence among the input features, which single remasking cannot address. To address this, we introduce multi-view feature remasking of each augmented graph $\mathcal{F}_{enc}^{\mathcal{G}_{i}^{prop}}$. Each remasked encoded feature is denoted by $\mathcal{F}_{enc_{j}}^{\mathcal{G}_{i}^{prop}}$. It acts as a regularizer for the decoder, making it robust against unexpected noises in input and helping to avoid overfitting. The final objective of the decoder is to reconstruct the actual node features \mathcal{F} from these masked encoded node features using the multi-view autoencoder loss described next.

Multi-View Autoencoder Loss: Given k augmentations of the input graph \mathcal{G}^{prop} represented as $\mathcal{G}_k^{prop}, \dots, \mathcal{G}_k^{prop}$, and m remarked decoded output for each augmented graph represented as

 $\mathcal{F}_{dec_1}^{\mathcal{G}_1^{prop}}, \dots, \mathcal{F}_{dec_m}^{\mathcal{G}_{ec_m}^{prop}}, \dots, \mathcal{F}_{dec_m}^{\mathcal{G}_k^{prop}}$, we define the multi-view reconstruction loss as

$$\mathcal{L}_{mrec} = \sum_{i=1}^{k} \sum_{j=1}^{m} \left| \left| \mathcal{F} - \mathcal{F}_{dec_j}^{\mathcal{G}_i^{prop}} \right| \right|_2^2 \tag{7}$$

To minimize the divergence across the views of the decoded features, we define the multi-view cosine similarity loss as

$$\mathcal{L}_{mcos} = \underset{\substack{\forall l, i, j; \text{ if } l = l' \text{ then } i \neq j \\ l \leq k, i \leq m, j \leq m}}{\mathcal{M}} \left(1 - \frac{\mathcal{F}_{dec}^{\mathcal{G}_{l}^{prop}}}{\left\| \mathcal{F}_{dec_{i}}^{\mathcal{G}_{l}^{prop}} \right\| \cdot \left\| \mathcal{F}_{dec_{j}}^{\mathcal{G}_{l'}^{prop}}} \right\|} \right)$$
(8)

Here, \mathcal{M} is the mean operation. Our final propagation loss is $\mathcal{L}_{prop} = \mathcal{L}_{mrec} + \mathcal{L}_{mcos}$.

3.3 Final Loss

We combine the AMR and Propagation loss as $\mathcal{L} = \mathcal{L}_{lgcl} + \mathcal{L}_{prop}$. We train our model using this loss, and the final features of our model are $\mathcal{H}^{\mathcal{G}^{amr}} \cdot \mathcal{H}^{\mathcal{G}^{prop}}$. These features are then used for misinformation classification.

4 Experiments and Results

Dataset Description and metrics: We perform experiments on the publicly available datasets Fake-NewsNet (Shu et al., 2020) in order to assess the effectiveness of the model. This repository contains two separate benchmark datasets, namely, PolitiFact and GossipCop. PolitiFact is dedicated to news coverage revolving around U.S. political affairs, while GossipCop has stories about Hollywood celebrities. These datasets also capture the broader social dynamics by including information about how news spreads through networks and the posting patterns of users. We evaluate our model using F1-score, and Accuracy (Acc). Comprehensive details of the datasets are provided in Table 2.

Table 2: Datasets Statistics

	# News	# True	# Fake	# Nodes	# Edges
PolitiFact	314	157	157	41054	40740
GossipCop	5464	2732	2732	314262	308798

Baselines: In our evaluation, we contrast our model with various state-of-the-art baselines, categorized into two groups. The first group utilizes only unsupervised methods (**TruthFinder** (Yin et al., 2008), **UFNDA** (Li et al., 2021), **UFD** (Yang et al., 2022), **GTUT** (Gangireddy et al., 2020), **(UMD)**² (Silva et al., 2024), **GraphMAE** (Hou et al., 2022), **GAMC** (Yin et al., 2024)), while the second incorporates supervised methods (**SAFE** (Zhou et al., 2020), **EANN** (Wang et al., 2018), **dEFEND** (Shu et al., 2019), **GACL** (Sun et al., 2022), **EA**²**N** (**BERT**) (Gupta et al., 2025)).

5 Results

We conducted a comparative analysis of our model against various baselines as mentioned above on the PolitiFact and GossipCop datasets. Results are shown in Table 3. Our model achieved the highest accuracy (0.919) and F1-score (0.918) among the unsupervised baselines. Compared to GAMC, the existing benchmark, our model outperforms it by a margin of 8.1% in accuracy and 8.7% in F1-score (on the absolute scale). Also, our model surpasses GTUT and $(UMD)^2$ by significant margins, $12 \sim 14\%$ in accuracy and $14 \sim 15\%$ in the F1-score, indicating a superior ability to differentiate between fake and real news. Similarly, our model significantly outperforms existing unsupervised baselines on the GossipCop dataset. It achieves the highest accuracy (0.968) and F1-score (0.966), outperforming GAMC, which attained an accuracy of 0.946 and an F1-score of 0.943. This represents a 2.2% improvement in accuracy and a 2.3% improvement in

the F1-score. This improvement can be attributed to the proposed model, which leverages a combination of self-supervised AMR semantic features and news propagation features from multi-view social context graph learning.

When we compare our model to supervised baselines on both PolitiFact and GossipCop datasets (Table 3), it consistently outperforms state-of-the-art approaches in terms of accuracy, while comparable results on F1 score are observed. On PolitiFact, our model achieves an accuracy of 0.919 and an F1-score of 0.933, surpassing EA²N with BERT (0.911 accuracy, 0.915 F1-score), GACL (0.867 accuracy, 0.866 F1-score), and EANN (0.804 accuracy, 0.798 F1-score). However, it shows comparative performance with dEFEND in F1-score. On GossipCop, our model outperforms all supervised baselines, achieving the highest accuracy (0.968) and F1-score (0.966). It notably surpasses GACL (0.907 accuracy, 0.905 F1-score) and EA²N (0.844 accuracy, 0.872 F1-score), as well as dEFEND, which lags significantly behind with 0.808 accuracy and 0.755 F1-score. These results highlight that while supervised models perform well, our self-supervised approach not only competes effectively on PolitiFact but outperforms all supervised baselines on GossipCop, demonstrating superior performance across datasets. Our self-supervised pipeline may yield stronger representations than shallow supervised models trained only on labels. One reason is that the datasets have known issues with label reliability. In such cases, supervised models can overfit to spurious correlations or unreliable labels and unsupervised models often rely on representation learning, which can be more robust to noise and generalize better in low-label regimes.

Table 3: Comparative study of our model w.r.t. different baselines on PolitiFact and GossipCop datasets.

Methods	Polit	iFact	GossipCop				
Methods	Acc	F1	Acc	F1			
Unsupervised Methods							
TruthFinder	0.581	0.573	0.668	0.669			
UFNDA	0.685	0.670	0.692	0.673			
$_{ m UFD}$	0.697	0.647	0.662	0.667			
GTUT	0.776	0.767	0.771	0.744			
$(\mathrm{UMD})^2$	0.802	0.761	0.792	0.783			
$\operatorname{GraphMAE}$	0.643	0.649	0.802	0.787			
GAMC	0.838	0.831	0.946	0.943			
Supervised Methods							
SAFE	0.793	0.775	0.832	0.811			
EANN	0.804	0.798	0.836	0.813			
dEFEND	0.904	0.928	0.808	0.755			
GACL	0.867	0.866	0.907	0.905			
$\mathrm{E}\mathrm{A}^2\mathrm{N}$	<u>0.911</u>	0.915	0.844	0.872			
Ours	0.919	0.918	0.968	0.966			
variance	± 0.019	± 0.020	± 0.015	± 0.015			

6 Ablation Study

Change in classification result with different values of λ : Figure 2 shows the change in classification accuracy of the method with the change in weightage to negative samples in Equation 4. It is evident that the accuracy improved initially with the value of λ and obtained the maximum result when $\lambda = 0.5$ for both datasets. With a further increase in λ , the accuracy decreases, indicating that our model overemphasizes negative samples compared to being close to positive samples, thus decreasing feature separability. Based on this, we set the value of λ to 0.5 in our experiments.

Change in classification result with training size: We conduct a classification experiment using a linear SVM with varying training sizes while keeping the test set fixed at 10%. The results shown in Table 4, clearly demonstrate that the classification accuracy improves as expected with larger training data. With just 10% of the training data, our model achieves superior performance on both the GossipCop (0.951 of

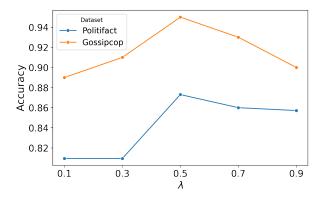


Figure 2: Change in classification result with different values of λ .

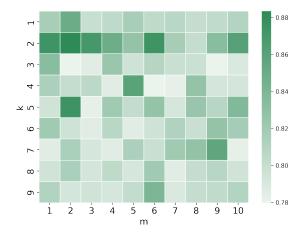


Figure 3: Change in accuracy with varying number of augmentation k and multi-view remarking m.

ours vs 0.946 of GAMC and 0.907 of GACL) and PolitiFact (0.875 of ours vs 0.838 of GAMC), highlighting its effectiveness in data-scarce scenarios. This showcases the robustness and generalization ability of the learned representations.

Change in results with varying number of augmentations k and multi-view remarkings m: We study the change in classification accuracy with different numbers of augmentations and remarkings for the PolitiFact dataset (Figure 3). We can infer from the figure that the best results are obtained when we set

Table 4: Results on different split sizes for PolitiFact and GossipCop datasets.

Train Size %	Polit	iFact	GossipCop		
Train Size 70	Acc	F1	Acc	F1	
10	0.875	0.867	0.951	0.951	
20	0.875	0.867	0.948	0.949	
30	0.875	0.867	0.951	0.951	
40	0.906	0.903	0.952	0.953	
50	0.906	0.903	0.952	0.953	
60	0.906	0.903	0.952	0.953	
70	0.906	0.909	0.952	0.953	
80	0.938	0.938	0.954	0.955	
90	0.938	0.941	0.956	0.957	

Model	Polit	iFact	GossipCop	
Wodel	Acc	F1	Acc	F1
Mistral (Zero-shot)	0.747	0.636	0.610	0.320
LLaMA (Zero-shot)	0.804	0.749	0.680	0.535
Only \mathcal{L}_{lgcl} + Mistral	0.822	0.830	0.934	0.932
Only \mathcal{L}_{lgcl} + LLaMA	0.841	0.828	0.948	0.949
Only \mathcal{L}_{prop}	0.846	0.845	0.946	0.945
$\mathcal{L}_{lgcl} + \mathcal{L}_{prop} + \text{Mistral}$	0.893	0.892	0.938	0.938
$\mathcal{L}_{lgcl} + \mathcal{L}_{prop} + \text{LLaMA}$	0.919	0.918	0.968	0.966

Table 5: Accuracy Score for different components of the model.

k=2 and $m\leq 6$. This shows that multi-view remarkings help the model achieve superior performance, but more than three remarkings do not bring considerable improvements.

Change in classification results with different components of our model: In Table 5, we show the importance of different components of our model. All the results shown here use 80% labelled data in the final linear SVM for training. As we can see from the table, \mathcal{L}_{lgcl} and \mathcal{L}_{prop} individually produce comparable results. But we get significant improvements in classification accuracy when we combine features generated using $\mathcal{L} = \mathcal{L}_{lgcl} + \mathcal{L}_{prop}$. We also compare the performance of our model with varying versions of the LLM, Mistral-7B and LLaMA-7B. Our model significantly improves the classification results using LLM-guided centroids and the proposed losses as compared to the independent LLMs. One must also note that there is a significant difference between the results from the two LLMs when used independently. But, when used with any component of our model, this difference reduces, thus showing the robustness of the extracted features by the proposed method.

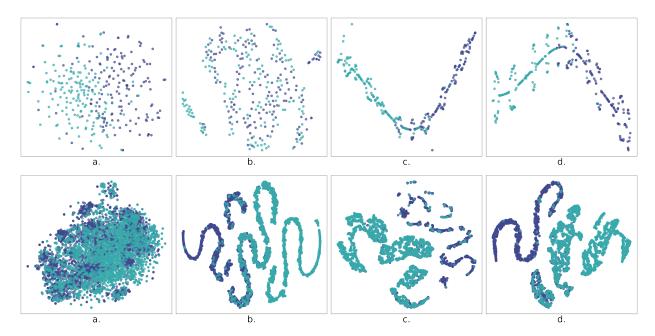


Figure 4: The TSNE plots showing the embeddings of PolitiFact (Row1) and GossipCop (Row2).

Qualitative results at different stages of our proposed pipeline In Figure 4 we show the feature separation between the real and fake news at different stages of our proposed pipeline. In the first row of the Figure we see the results of PolitiFact dataset and the second row we show the results of the GossipCop dataset. The first column of each row shows the TSNE embedding of the initial features. The second column shows the TSNE plot of the original features after a single fully connected linear layer (MLP). The third column shows the TSNE plot of the features obtained after the self-supervised AMR graph learning ($\mathcal{H}^{\mathcal{G}^{amr}}$) phase trained with a linear layer. The last columns shows the TSNE plot of the final concatenated features after self-supervised AMR graph learning and multi-view propagation graph learning ($\mathcal{H}^{\mathcal{G}^{amr}}.\mathcal{H}^{\mathcal{G}^{prop}}$) with a linear layer. In all the cases we train the MLP with 80% labelled data.

To quantify the clustering quality, we compute the silhouette score at each stage. For the PolitiFact dataset, the silhouette scores are 0.33, 0.54, 0.62, and 0.64, respectively, indicating progressively better separation between real and fake news as the pipeline advances. Similarly, for the GossipCop dataset, the silhouette scores are 0.16, 0.34, 0.38, and 0.40, again demonstrating consistent improvement. These quantitative results further support the visual evidence, confirming that our model increasingly enhances feature discriminability at each stage of the pipeline.

7 Implementation Details

In order to generate the AMR graph, we have used a pretrained STOG model (Zhang et al., 2019). For LGCL, we use $\alpha=0.5$ and in order to integrate the evidence in the AMR graph, we use the same parameters described in Gupta et al. (2025). For social context and propagation graph learning we use 2 encoder layers and 1 decoder layer. For multi-view remasking, we select k=2 and m=2. We selected Support Vector Machine (SVM) as the classifier in the downstream task and reported the results from 80 % of the training data with 5-fold cross-validation. We have trained our model on RTX A5000 NVIDIA GPU with 24 GB GPU memory. The training of AMR took 1 hour for PolitiFact and took 3 hours for the GossipCop dataset with 50 epochs. Multi-view masked graph learning took 5 mins for the PolitiFact dataset and 15 minutes for the GossipCop dataset.

8 Conclusion

This study presents a novel self-supervised approach for misinformation detection. The LLM-guided contrastive self-supervised AMR learning framework captures complex semantic relationships in text. This method enhances feature separation between real and fake news by leveraging an LLM-guided negative sampler. Additionally, we introduce a multi-view graph-masked autoencoder that integrates social context and news propagation patterns for more robust detection. Through extensive experiments, the proposed method is found to produce state-of-the-art performance.

References

Laura Banarescu, Claire Bonial, Shu Cai, Madalina Georgescu, Kira Griffitt, Ulf Hermjakob, Kev Knight, Philipp Koehn, Martha Palmer, and Nathan Schneider. Abstract Meaning Representation for sembanking. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pp. 178–186, Sofia, Bulgaria, August 2013.

Deng Cai and Wai Lam. Graph transformer for graph-to-sequence learning. In AAAI, pp. 7464–7471. AAAI Press, 2020.

Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *EMNLP*, pp. 1724–1734, Doha, Qatar, October 2014. ACL.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pp. 4171-4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL https://aclanthology.org/N19-1423/.

Hernawan Dewatana and Siti Ummu Adillah. The effectiveness of criminal eradication on hoax information and fake news. Law Development Journal, 3(3):513–520, 2021.

- Yingtong Dou, Kai Shu, Congying Xia, Philip S. Yu, and Lichao Sun. User preference-aware fake news detection. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '21, pp. 2051–2055, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450380379. doi: 10.1145/3404835.3462990. URL https://doi.org/10.1145/3404835.3462990.
- Yaqian Dun, Kefei Tu, Chen Chen, Chunyan Hou, and Xiaojie Yuan. Kan: Knowledge-aware attention network for fake news detection. AAAI, 35(1):81–89, May 2021.
- Song Feng, Ritwik Banerjee, and Yejin Choi. Syntactic stylometry for deception detection. In *ACL (Volume 2: Short Papers)*, pp. 171–175, Jeju Island, Korea, July 2012. ACL.
- Siva Charan Reddy Gangireddy, Deepak P, Cheng Long, and Tanmoy Chakraborty. Unsupervised fake news detection: A graph-based approach. In *Proceedings of the 31st ACM Conference on Hypertext and Social Media*, HT '20, pp. 75–83, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450370981. doi: 10.1145/3372923.3404783. URL https://doi.org/10.1145/3372923.3404783.
- Bilal Ghanem, Simone Paolo Ponzetto, Paolo Rosso, and Francisco Rangel. Fakeflow: Fake news detection by modeling the flow of affective information. In 16th EACL, 2021.
- Shubham Gupta, Narendra Yadav, Suman Kundu, and Sainathreddy Sankepally. Fakedamr: Fake news detection using abstract meaning representation network. In *International Conference on Complex Networks and Their Applications*, pp. 308–319. Springer, 2023.
- Shubham Gupta, Abhishek Rajora, and Suman Kundu. Ea2n: Evidence-based amr attention network for fake news detection. *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–12, 2025. doi: 10.1109/TKDE.2025.3529707.
- Zhenyu Hou, Xiao Liu, Yukuo Cen, Yuxiao Dong, Hongxia Yang, Chunjie Wang, and Jie Tang. Graphmae: Self-supervised masked graph autoencoders. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '22, pp. 594–604, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450393850. doi: 10.1145/3534678.3539321. URL https://doi.org/10.1145/3534678.3539321.
- Zhenyu Hou, Yufei He, Yukuo Cen, Xiao Liu, Yuxiao Dong, Evgeny Kharlamov, and Jie Tang. Graphmae2: A decoding-enhanced masked self-supervised graph learner. In *Proceedings of the ACM Web Conference 2023*, WWW '23, pp. 737–746, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9781450394161. doi: 10.1145/3543507.3583379. URL https://doi.org/10.1145/3543507.3583379.
- Linmei Hu, Tianchi Yang, Luhao Zhang, Wanjun Zhong, Duyu Tang, Chuan Shi, Nan Duan, and Ming Zhou. Compare to the knowledge: Graph neural fake news detection with external knowledge. In *ACL-IJCNLP* (Volume 1: Long Papers), pp. 754–763, Online, August 2021. ACL.
- Thomas N. Kipf and Max Welling. Variational graph auto-encoders, 2016. URL https://arxiv.org/abs/1611.07308.
- Dun Li, Haimei Guo, Zhenfei Wang, and Zhiyun Zheng. Unsupervised fake news detection based on autoencoder. *IEEE Access*, 9:29356–29365, 2021. doi: 10.1109/ACCESS.2021.3058809.
- Shaohua Li, Weimin Li, Alex Munyole Luvembe, and Weiqin Tong. Graph contrastive learning with feature augmentation for rumor detection. *IEEE Transactions on Computational Social Systems*, 11(4):5158–5167, 2024. doi: 10.1109/TCSS.2023.3269303.
- Zewen Li, Fan Liu, Wenjie Yang, Shouheng Peng, and Jun Zhou. A survey of convolutional neural networks: Analysis, applications, and prospects. *IEEE Transactions on Neural Networks and Learning Systems*, 33 (12):6999–7019, 2022. doi: 10.1109/TNNLS.2021.3084827.
- Yang Liu and Yi-Fang Wu. Early detection of fake news on social media through propagation path classification with recurrent and convolutional networks. AAAI, 32(1), Apr. 2018.

- Yunfei Long, Q Lu, Rong Xiang, Minglei Li, and Chu-Ren Huang. Fake news detection through multiperspective speaker profiles. In *IJCNLP* (*Volume 2: Short Papers*), pp. 252–256, Taipei, Taiwan, November 2017. Asian Federation of Natural Language Processing.
- Jing Ma, Wei Gao, Prasenjit Mitra, Sejeong Kwon, Bernard J. Jansen, Kam-Fai Wong, and Meeyoung Cha. Detecting rumors from microblogs with recurrent neural networks. In *IJCAI*, IJCAI'16, pp. 3818–3824. AAAI Press, 2016a. ISBN 9781577357704.
- Jing Ma, Wei Gao, Prasenjit Mitra, Sejeong Kwon, Bernard J. Jansen, Kam-Fai Wong, and Meeyoung Cha. Detecting rumors from microblogs with recurrent neural networks. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, IJCAI'16, pp. 3818–3824. AAAI Press, 2016b. ISBN 9781577357704.
- Erxue Min, Yu Rong, Yatao Bian, Tingyang Xu, Peilin Zhao, Junzhou Huang, and Sophia Ananiadou. Divide-and-conquer: Post-user interaction network for fake news detection on social media. In *Proceedings of the ACM Web Conference 2022*, WWW '22, pp. 1148–1158, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450390965. doi: 10.1145/3485447.3512163. URL https://doi.org/10.1145/3485447.3512163.
- Kashyap Popat, Subhabrata Mukherjee, Jannik Strötgen, and Gerhard Weikum. Where the truth lies: Explaining the credibility of emerging claims on the web and social media. WWW '17 Companion, pp. 1003–1012, Republic and Canton of Geneva, CHE, 2017. International World Wide Web Conferences Steering Committee. ISBN 9781450349147.
- Kai Shu, Limeng Cui, Suhang Wang, Dongwon Lee, and Huan Liu. defend: Explainable fake news detection. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '19, pp. 395–405, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450362016. doi: 10.1145/3292500.3330935. URL https://doi.org/10.1145/3292500.3330935.
- Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dongwon Lee, and Huan Liu. Fakenewsnet: A data repository with news content, social context, and spatiotemporal information for studying fake news on social media. *Biq Data*, 8(3):171–188, 2020.
- Amila Silva, Ling Luo, Shanika Karunasekera, and Christopher Leckie. Unsupervised Domain-Agnostic Fake News Detection Using Multi-Modal Weak Signals. *IEEE Transactions on Knowledge & Data Engineering*, 36(11):7283-7295, November 2024. ISSN 1558-2191. doi: 10.1109/TKDE.2024.3392788. URL https://doi.ieeecomputersociety.org/10.1109/TKDE.2024.3392788.
- Tiening Sun, Zhong Qian, Sujun Dong, Peifeng Li, and Qiaoming Zhu. Rumor detection on social media with graph adversarial contrastive learning. In *Proceedings of the ACM Web Conference 2022*, WWW '22, pp. 2789–2797, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450390965. doi: 10.1145/3485447.3511999. URL https://doi.org/10.1145/3485447.3511999.
- Yaqing Wang, Fenglong Ma, Zhiwei Jin, Ye Yuan, Guangxu Xun, Kishlay Jha, Lu Su, and Jing Gao. Eann: Event adversarial neural networks for multi-modal fake news detection. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '18, pp. 849–857, New York, NY, USA, 2018. Association for Computing Machinery. ISBN 9781450355520. doi: 10.1145/3219819.3219903. URL https://doi.org/10.1145/3219819.3219903.
- Lirong Wu, Haitao Lin, Cheng Tan, Zhangyang Gao, and Stan Z. Li. Self-supervised learning on graphs: Contrastive, generative, or predictive. 35(4):4216–4235, April 2023. ISSN 1041-4347. doi: 10.1109/TKDE. 2021.3131584. URL https://doi.org/10.1109/TKDE.2021.3131584.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=ryGs6iA5Km.

- Ruichao Yang, Xiting Wang, Yiqiao Jin, Chaozhuo Li, Jianxun Lian, and Xing Xie. Reinforcement subgraph reasoning for fake news detection. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '22, pp. 2253–2262, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450393850. doi: 10.1145/3534678.3539277. URL https://doi.org/10.1145/3534678.3539277.
- Huaxiu Yao, Ying-xin Wu, Maruan Al-Shedivat, and Eric Xing. Knowledge-aware meta-learning for low-resource text classification. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 1814–1821, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
- Shu Yin, Peican Zhu, Lianwei Wu, Chao Gao, and Zhen Wang. Gamc: An unsupervised method for fake news detection using graph autoencoder with masking. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(1):347–355, Mar. 2024. doi: 10.1609/aaai.v38i1.27788. URL https://ojs.aaai.org/index.php/AAAI/article/view/27788.
- Xiaoxin Yin, Jiawei Han, and Philip S. Yu. Truth discovery with multiple conflicting information providers on the web. *IEEE Transactions on Knowledge and Data Engineering*, 20(6):796–808, 2008. doi: 10.1109/TKDE.2007.190745.
- Chunyuan Yuan, Qianwen Ma, Wei Zhou, Jizhong Han, and Songlin Hu. Jointly Embedding the Local and Global Relations of Heterogeneous Graph for Rumor Detection. In 2019 IEEE International Conference on Data Mining (ICDM), pp. 796–805, Los Alamitos, CA, USA, November 2019. IEEE Computer Society. doi: 10.1109/ICDM.2019.00090. URL https://doi.ieeecomputersociety.org/10.1109/ICDM.2019.00090.
- Sheng Zhang, Xutai Ma, Kev Duh, and Benjamin Van Durme. AMR parsing as sequence-to-graph transduction. In ACL, pp. 80–94, Florence, Italy, July 2019. ACL.
- Yizhou Zhang, Loc Trinh, Defu Cao, Zijun Cui, and Yan Liu. Detecting out-of-context multimodal misinformation with interpretable neural-symbolic model, 2023.
- Xinyi Zhou, Jindi Wu, and Reza Zafarani. Safe: Similarity-aware multi-modal fake news detection. In Hady W. Lauw, Raymond Chi-Wing Wong, Alexandros Ntoulas, Ee-Peng Lim, See-Kiong Ng, and Sinno Jialin Pan (eds.), *Advances in Knowledge Discovery and Data Mining*, pp. 354–367, Cham, 2020. Springer International Publishing. ISBN 978-3-030-47436-2.