

# MGM: Global Understanding of Audience Overlap Graphs for Predicting the Factuality and the Bias of News Media

Anonymous ACL submission

## Abstract

In the current era of rapidly growing digital data, evaluating the political bias and factuality of news media outlets has become more important for seeking reliable information online. In this work, we study the classification problem of profiling the news media from the lens of political bias and factuality. Traditional profiling methods, such as Graph Neural Networks (GNNs) and Pre-trained Language Models (PLMs), have shown promising results, but they face notable challenges. GNNs often struggle with media graphs containing disconnected components and insufficient labels, while PLMs focus solely on textual features, causing them to overlook the complex relationships between entities. To address these limitations, we propose MediaGraphMind (MGM), an effective solution within a variational Expectation-Maximization (EM) framework. Instead of relying on limited neighboring nodes, MGM leverages features, structural patterns, and label information from globally similar nodes. Such a framework not only enables GNNs to capture long-range dependencies for learning expressive node representations but also enhances PLMs by integrating structural information and therefore improving the performance of both models. Our extensive experiments on a standard dataset demonstrate the effectiveness of our proposed framework.

## 1 Introduction

The rise of the Internet has offered many opportunities to publish information and to express opinions (Nakov et al., 2021). At the same time, it has accelerated the spread of misinformation and disinformation online (Fairbanks et al., 2018). Initially, efforts focused on verifying individual claims, but it soon became evident that assessing the factuality of the news sources themselves was equally important (Baly et al., 2020b).

Early studies on automatic media profiling relied solely on text characteristics (Battaglia et al., 2018;

Pérez-Rosas et al., 2017), which has proven particularly challenging (Baly et al., 2018, 2020b). The complexity is heightened when the text features contain indeterminate noise, leading to classification errors. Moreover, traditional methods struggle to capture the intricate relationships between entities, such as articles, audiences, and media outlets. Graph Neural Networks (GNNs) have emerged as an effective framework to capture and to model these relationships, reducing the dependence on individual text features and improving media profiling (Panayotov et al., 2022; Mehta et al., 2022).

Despite their advantages, GNNs also face challenges due to the graph being disconnected and due to label sparsity. *Disconnected components* hinder GNNs from capturing long-range dependencies, thus limiting their ability to learn expressive node representations (Bodnar et al., 2021; Dai et al., 2021, 2022). To address these issues, various memory-based GNNs have been proposed (Kang, 2021; Zhang et al., 2022; Ma et al., 2022), which capture long-range dependencies using external memory modules that store *global information*, i.e., features and structural patterns throughout the graph. However, these approaches require substantial memory to store the embeddings of all nodes.

To address these challenges, we propose MGM, a method with a variational Expectation-Maximization (EM) framework that enhances existing GNNs to capture and leverage global information in media graphs. MGM captures both local and global patterns, node features, and labels from global similar nodes for improved performance. Unlike Graph Attention Networks (Veličković et al., 2018), which focus on similar nodes in local neighbors, MGM uses an external memory module to store precomputed node representations, thereby reducing compute costs (Fey et al., 2021) and enabling efficient retrieval of node embeddings. Moreover, MGM reduces memory requirements by focusing on a small set of candidate nodes, guided

by a distribution over the training nodes sampled from a Dirichlet prior (He et al., 2020).

The results show that MGM significantly improves the performance of the baseline GNNs, achieving a 10% improvement in all evaluation measures on the Media Bias/Fact Check (MBFC)<sup>1</sup> data featured in the ACL-2020 (Baly et al., 2020c) and the EMNLP-2018 (Baly et al., 2018) datasets. Despite the lack of rich features in the graph, we enhance the dataset by scraping *Articles* and *Wikipedia* for ACL-2020. Pre-trained language models (PLMs) such as BERT (Devlin et al., 2018), RoBERTa (Liu et al., 2019), DistilBERT (Sanh et al., 2019), and DeBERTaV3 (He et al., 2021) are fine-tuned to predict political bias and factuality. Where media data are inaccessible, MGM’s representation-based probabilities fill the gap. Moreover, integrating MGM’s probabilities with LM’s enhances the performance for both tasks. Our contributions can be summarised as follows:

- We introduce MGM, an effective approach that leverages global information to enhance the expressiveness of existing GNNs for reliable news media profiling.
- In MGM, we use external memory to efficiently retrieve globally similar nodes and reduce memory requirements by learning a sparse distribution with a Dirichlet prior.
- We illustrate that MGM consistently outperforms vanilla GNNs for the detection of factuality and political bias across all baselines.
- We validated that integrating the MGM features with the PLMs enhances performance and yields state-of-the-art results.

## 2 Related Work

### 2.1 Factuality and Political Bias of Media

Early research on news media profiling focused on textual content analysis (Afroz et al., 2012; Battaglia et al., 2018; Pérez-Rosas et al., 2017; Conroy et al., 2015). To improve the performance, subsequent research added contextual information (Baly et al., 2020c; Hounsel et al., 2020; Castelo et al., 2019; Fairbanks et al., 2018), including the nuances of multimedia production (Huh et al., 2018), the associated infrastructure (Hounsel et al., 2020), and the social context (Baly et al.,

2020c). Guo et al. (2022) used BERT (Devlin et al., 2019) to model the linguistic political bias in news articles. Fan et al. (2019) used annotated media from Budak et al. (2016), analyzing articles for political bias using distant supervision. Various methods measured political bias, including analyzing Twitter interactions (An et al., 2011, 2012; Stefanov et al., 2020), often using small English-only datasets (Kulkarni et al., 2018; Potthast et al., 2018; Kiesel et al., 2019; Baly et al., 2020a; Da San Martino et al., 2023; Nakov et al., 2023a,b; Barrón-Cedeño et al., 2023a,b; Azizov et al., 2023; Chen et al., 2018; Fan et al., 2019; Spinde et al., 2022).

The veracity of the news media has been explored using PLMs to estimate the reliability of the source, correlated with the ratings of human experts (Yang and Menczer, 2023). Mehta and Goldwasser (2023) introduced a framework that combines graph-based models, PLMs, and human experience to profile news media, effectively identifying fake news with minimal human input. Early work estimated the reliability of the source based on the stance toward true/false claims using English datasets (Mukherjee and Weikum, 2015; Popat et al., 2017). Recent approaches, such as Baly et al. (2020c), used gold labels and various English sources as features to profile media with PLMs. Although the features in the aforementioned studies are obtained from diverse sources, they neglect the inherent relationships between the media.

To bridge this gap, graphs emerged as a comprehensive and effective framework for representation learning (Mehta et al., 2022). Panayotov et al. (2022) constructed a graph based on the principle of homophily, suggesting that similar media sources attract similar audiences. The framework leveraged the audience overlap of media outlets to build a huge graph that models the interactions between media, and to learn expressive representation for the nodes using GNNs. However, media graphs are characterized by disconnected components and scarce labels. To overcome these limitations, we propose MGM to capture long-range dependencies, and to improve the effective use of information across the entire graph.

### 2.2 Graph Neural Networks

The current *de facto* design of GNNs follows the message-passing framework (Yang et al., 2022), where they learn node representations by aggregating information from local neighbors. However, media graphs suffer from challenges such

<sup>1</sup>[www.mediabiasfactcheck.com](http://www.mediabiasfactcheck.com)

as multiple disconnected components and limited labels, which makes it difficult for GNNs to capture long-range dependencies and to learn informative node representations (Cui et al., 2020; Sun et al., 2020; Bodnar et al., 2021; Dai et al., 2021, 2022). Recent efforts on the integration of external memory modules to store embeddings of all nodes allow GNNs to capture long-range dependencies across graphs (Kang, 2021; Zhang et al., 2022; Ma et al., 2022). However, these methods typically require storing embeddings for all nodes in the graph, which results in high memory costs and low efficiency during testing. Unlike previous approaches, MGM focuses on a small set of candidate nodes, which are more likely to be selected as global similar nodes based on a Dirichlet prior distribution applied to the training nodes (Sethuraman, 1994).

### 3 Problem Formulation

We instantiate the graph learning task with the node classification problem in a semi-supervised setting (Kipf and Welling, 2016; Veličković et al., 2018), where labels are only available for a small subset of the nodes. Let  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}, \mathbf{X}, \mathbf{Y}^l\}$  represent a partially-labeled graph, where  $\mathcal{V} = \{v_i\}_{i=1}^N$  is a set of nodes,  $\mathcal{E}$  is a set of edges, and  $N$  is the total number of nodes. The node features are denoted as  $\mathbf{X} \in \mathbb{R}^{N \times F}$ , where  $F$  is the feature dimension. Since most nodes are unlabeled,  $\mathcal{V}$  can be divided into labeled nodes  $\mathcal{V}^l$  with labels  $\mathbf{Y}^l$ , and unlabeled nodes  $\mathcal{V}^u$ . The labels  $\mathbf{Y}^l \in \mathbb{R}^{N_l \times C}$  are in a one-hot form, where  $N_l$  and  $C$  represent the number of labeled nodes and the number of classes, respectively. The goal of semi-supervised learning is to learn the model parameters  $\theta$  by maximizing the marginal distribution of the overall labeled nodes, i.e.,  $p_\theta(\mathbf{Y}^l | \mathbf{X}, \mathcal{E}) = \prod_{n \in \mathcal{V}^l} p_\theta(\mathbf{y}_n | \mathbf{X}, \mathcal{E})$  on the training graph.

### 4 Methodology

In this section, we introduce our proposed framework that aims to enhance the performance of GNNs and PLMs for factuality and political bias. GNNs struggle with media graphs containing disconnected components and insufficient labels, while PLMs only use textual features, neglecting the intricate relationships between entities. Thus, we aim to develop a framework that not only enhances the performance of GNNs, but also integrates structural information into the PLMs.

Toward this goal, we propose our novel MGM

framework, which uses features, structural patterns, and label information of global similar nodes to enhance the performance of GNNs for factuality and political bias classification. Moreover, we integrate MGM with PLMs to further improve the performance. With such a framework, GNN can effectively overcome the challenges in media graphs. Meanwhile, PLMs can benefit from integrated structural information, leading to enhanced performance for our tasks.

#### 4.1 The MGM Framework

Following (Qu et al., 2019, 2021), we adopt a probabilistic framework for node classification, treating node representations  $\mathbf{Z}$  as latent variables determined by a GNN. To improve the performance of the model, we propose to augment the GNNs with information about *global similar nodes*, i.e., nodes in the entire graph that have similar node features and local geometric structures. Specifically, we denote the set of global similar nodes of node  $n$  as  $\mathbf{t}_n \in \{0, 1\}^{N_l}$ , where  $\mathbf{t}_{nm} = 1$  indicates that node  $m$  is a global node similar to  $n$ . Similarly to node representations, we also regard the similar node indicator  $\mathbf{t}_n$  as a latent variable. Therefore, the joint probability distribution of global information-enhanced method can be factorized as follows:

$$\begin{aligned} p_\theta(\mathbf{Y}^l, \mathbf{T}, \mathbf{Z} | \mathbf{X}, \mathcal{E}) &= \\ &= p_\theta(\mathbf{Z} | \mathbf{X}, \mathcal{E}) p_\theta(\mathbf{T} | \mathbf{Z}) p_\theta(\mathbf{Y}^l | \mathbf{Z}, \mathbf{T}), \end{aligned} \quad (1)$$

where  $\mathbf{T} = [\mathbf{t}_n]_{n \in \mathcal{V}^l}^\top$  are the global similar nodes of all nodes.

However, finding global similar nodes with node representations requires computing representations for all nodes, which is expensive in terms of space and time (Fey et al., 2021). To alleviate this, we propose to store the embeddings of the labeled nodes in the memory and to use them to find global similar nodes. Consequently, the distribution of  $\mathbf{T}$  can be replaced by  $p_\theta(\mathbf{T} | \hat{\mathbf{Z}})$ , where  $\hat{\mathbf{Z}}$  is the embeddings of the labeled nodes in the memory, i.e.,  $\hat{\mathbf{Y}}^l$ . In this case, we can directly retrieve the representation from memory without computing representations for all nodes, thus making it more efficient to obtain the distribution of global similar nodes for both training and prediction.

To reduce the memory size, we select global similar nodes from a small set of candidate nodes, which are a subset of the training nodes. As a result, only the embeddings of these candidate nodes are stored in memory for prediction. To achieve

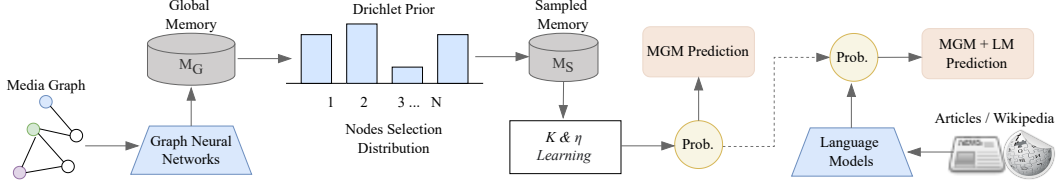


Figure 1: Key components of our proposed approach. GNNs store the representation of the media graphs in an external global memory ( $M_g$ ). A Dirichlet prior is used to select the distribution of sparse candidate nodes, which are stored in the sampled memory ( $M_s$ ). The parameters  $K$  and  $\eta$  control the number of candidate nodes and their influence, balancing local and global information. Since PLMs miss some of the media representation, they leverage MGM representation-based probabilities for the classification task. Detailed pipeline of the integration MGM with PLMs could be seen in Figure 3 (Appendix E).

this, we assume that  $p_\theta(\mathbf{T} \mid \hat{\mathbf{Z}})$  is a sparse distribution, concentrated on a few candidate nodes. Since the candidate set is not known, we introduce a latent variable  $\omega$  for each node  $n$ , where  $\omega_i \in [0, 1]$ , s.t.  $\sum_{i=1}^{N_i} \omega_i = 1$ . Here,  $\omega_i$  represents the probability that the  $i$ -th node in the labeled node is a candidate node. Inspired by (He et al., 2020), we introduce a prior over  $\omega$ , i.e.,  $p_\alpha(\omega)$  with parameter  $\alpha$ . This prior is designed to encourage a sparse distribution over  $\omega$ . Therefore, the joint distribution of the method is now defined as follows:

$$p_\theta(\mathbf{Y}^l, \mathbf{T}, \mathbf{Z}, \omega \mid \mathbf{X}, \mathcal{E}, \hat{\mathbf{Z}}) = p_\alpha(\omega) \quad (2)$$

$$p_\theta(\mathbf{Z} \mid \mathbf{X}, \mathcal{E}) p_\theta(\mathbf{T} \mid \omega, \hat{\mathbf{Z}}) p_\theta(\mathbf{Y}^l \mid \mathbf{T}, \mathbf{Z}).$$

Next, we introduce the parameterization of our probabilistic framework.

**Prior distribution over  $\omega$ .** We use the Dirichlet distribution as the prior distribution over  $\omega$ , i.e.,  $p_\alpha(\omega) \propto \prod_{i=1}^N \omega_i^{\alpha_i - 1}$ , where  $\alpha_i$  is the concentration parameter of the distribution. The concentration parameter  $\alpha$  is a positive value and a smaller value of  $\alpha$  prefers a sparser distribution over  $\omega$  (He et al., 2020). In our experiments, we set  $\alpha < 1$  to encourage the sparse nodes distribution.

**Prior distribution over node representations  $\mathbf{Z}$ .** We model the prior distribution over node representations as Gaussian distributions (Bojchevski and Günnemann, 2018), which are obtained with GNNs due to their effectiveness in graph learning tasks. Therefore, the prior distribution over  $\mathbf{Z}$  is defined as follows:

$$p_\theta(\mathbf{Z} \mid \mathbf{X}, \mathcal{E}) = \mathcal{N}(\mathbf{Z} \mid \text{GNN}_\theta(\mathbf{X}, \mathcal{E}), \sigma_1^2 \mathbf{I}), \quad (3)$$

where  $\sigma_1^2$  is the learned variance of the prior and  $\text{GNN}_\theta$  is an  $L$ -layer GNN with parameter  $\theta$ .

**Prior distribution over  $\mathbf{T}$ .** To obtain global similar nodes of node  $n$ , we define a prior distribution

over  $\mathbf{T}$  as follows:

$$p_\theta(\mathbf{T} \mid \omega, \hat{\mathbf{Z}}) = \text{Mul}(\mathbf{T} \mid K, f_\theta(\omega, \hat{\mathbf{Z}})), \quad (4)$$

where  $\text{Mul}(\cdot)$  represents the multinomial distribution,  $K$  denotes the predefined number of global similar nodes, and  $f_\theta$  is designed as a parameterized function that outputs the parameters of the multinomial distribution.

**Prediction of label  $\mathbf{Y}$ .** Finally, we use the node representation  $\mathbf{Z}$  and information from global similar nodes to predict the label. Specifically, we leverage the labels of global similar nodes and first predict the label based on its representation:

$$p_\theta(\mathbf{Y} \mid \mathbf{Z}) = \text{Cat}(\mathbf{Y} \mid \mathbf{Z}), \quad (5)$$

where  $p_\theta(\mathbf{Y} \mid \mathbf{Z})$  is formulated as a categorical distribution. Then, we predict the label using the labels of global similar nodes:

$$p_\theta(\mathbf{Y} \mid \mathbf{T}) \propto \sum_{\mathbf{N}, \mathbf{M} \in \hat{\mathcal{V}}^l} \mathbf{T}_{\mathbf{N}\mathbf{M}} \cdot \mathbf{Y}_{\mathbf{M}}, \quad (6)$$

where  $\mathbf{T}_{\mathbf{N}\mathbf{M}}$  represents the indices of global similar nodes for the predicted nodes set  $\mathbf{N}$ . Additionally,  $\mathbf{Y}_{\mathbf{M}}$  denotes the one-hot labels of the nodes in  $\mathbf{M}$ , where  $\mathbf{M}$  is the set of global similar nodes. Finally, the predicted label distribution is defined as follows:

$$p_\theta(\mathbf{Y} \mid \mathbf{Z}, \mathbf{T}) = \eta p_\theta(\mathbf{Y} \mid \mathbf{Z}) + (1 - \eta) p_\theta(\mathbf{Y} \mid \mathbf{T}), \quad (7)$$

where  $\eta \in [0, 1]$  is a trade-off hyper-parameter. When  $\eta = 1$ , our model only uses local representations of nodes for prediction, which degrades to vanilla GNNs. In contrast, when  $0 < \eta < 1$ , our model predicts the labels of the nodes using information from both local neighbors and global similar nodes.



## 4.2 Training Process of MGM

Next, we explain how to learn the model parameters  $\theta$  based on the graph. Ideally, the marginal likelihood should be optimized during training:

$$p_\theta(\mathbf{Y}^l | \mathbf{X}, \mathcal{E}, \hat{\mathbf{Z}}) = \int_{\omega} \int_{\mathbf{Z}} \sum_{\mathbf{T}} p_\theta(\mathbf{Y}^l, \mathbf{T}, \mathbf{Z}, \omega | \mathbf{X}, \mathcal{E}, \hat{\mathbf{Z}}) d\mathbf{Z} d\omega. \quad (8)$$

However, the computation of maximizing the marginal likelihood is intractable due to the marginalization of latent variables. Therefore, we develop a variational Expectation-Maximization (EM) algorithm (Qu et al., 2019) to optimize its evidence lower bound (ELBO) instead:

$$\begin{aligned} \mathcal{L}_{\text{ELBO}}(\mathbf{Y}^l; \theta, \phi, \alpha, \lambda) = & -D_{\text{KL}}(q_\lambda(\omega) || p_\alpha(\omega)) \\ & - D_{\text{KL}}[q_\phi(\mathbf{Z} | \mathbf{T}, \mathbf{Y}^l) || p_\theta(\mathbf{Z} | \mathbf{X}, \mathcal{E})] \\ & - D_{\text{KL}}[q_\phi(\mathbf{T} | \mathbf{Y}^l) || p_\theta(\mathbf{T} | \omega, \hat{\mathbf{Z}})] \\ & + \mathbb{E}_{q_\phi(\mathbf{T} | \mathbf{Y}^l) q_\phi(\mathbf{Z} | \mathbf{T}, \mathbf{Y}^l)} [\log p_\theta(\mathbf{Y}^l | \mathbf{T}, \mathbf{Z})], \end{aligned} \quad (9)$$

where  $D_{\text{KL}}[\cdot || \cdot]$  is the Kullback-Leibler (KL) divergence,  $q$  represents the variational distribution to approximate the model posterior distribution and adheres to the following factorization form:<sup>2</sup>

$$q_\lambda(\omega) q_\phi(\mathbf{T}, \mathbf{Z}, \omega | \mathbf{Y}^l) q_\phi(\mathbf{T} | \mathbf{Y}^l) q_\phi(\mathbf{Z} | \mathbf{T}, \mathbf{Y}^l),$$

where  $\phi$  and  $\lambda$  are variational parameters.

Note that we use the mean-field assumption to approximate the posterior of  $\omega$  to simplify the variational distributions. For computational convenience, we assume that the variational distributions of these latent variables have the same distribution form as their prior distributions. Therefore, we define the variational distributions of  $\omega$ ,  $\mathbf{T}$  and  $\mathbf{Z}$  to be Dirichlet, multinomial, and Gaussian distributions, respectively.

Note that the KL divergence in Equation (9) has a closed-form solution, and we approximate the expectation using a Monte Carlo method by sampling from the variational distributions. In variational EM, the variational parameters  $\phi$  and the model parameters  $\theta$  are learned alternately. In the E-step, we fix  $\theta$  and update  $\phi$  by minimizing the KL divergence to approximate the true posteriors. In the M-step, we fix  $\phi$  and update  $\theta$  by maximizing the expected log-likelihood.

<sup>2</sup>We omit the dependence of variational distributions on node features  $\mathbf{X}$ , edges  $\mathcal{E}$  and memory  $\hat{\mathbf{Z}}$  for brevity.

## 4.3 Prediction Process of MGM

After training, we expect to obtain a sparse distribution  $q_\lambda(\omega)$ , allowing us to select a subset of the candidate nodes. In this case, we can select candidate nodes over a certain probability threshold, thus reducing the memory size and improving the efficiency for prediction. Specifically, we calculate the expected value of  $q_\lambda(\omega)$  for each node  $i$ , which is given by  $\mathbb{E}_{q_\lambda(\omega)}[\omega_i] = \lambda_i / \sum_{j=1}^{N_l} \lambda_j$ , and then we select the top- $M$  nodes that occupy 90% of the probability mass as candidate nodes.

We then leverage the embeddings of the memorized candidate nodes  $\hat{\mathbf{Z}}_\omega$  and  $p_\theta$  to predict the labels of the test nodes  $\tilde{n}$  based on Equation (7).

## 4.4 Enhancing PLMs Predictions with MGM

Next, we demonstrate how MGM improves the performance of PLMs by incorporating information from global similar nodes. Given textual features  $\mathbf{S}$ , such as those from *Articles* and *Wikipedia* pages for the media outlet, we first fine-tune the PLMs using the cross-entropy loss. Then, we concatenate the predicted label distribution from the PLMs with MGM to obtain the final label distribution:

$$\begin{aligned} p_{\psi, \theta}(\mathbf{Y} | \mathbf{S}, \mathbf{Z}, \mathbf{T}) = & \quad (10) \\ = & \text{Softmax}(\oplus(p_\psi(\mathbf{Y} | \mathbf{S}), p_\theta(\mathbf{Y} | \mathbf{Z}, \mathbf{T}))\mathbf{W} + \mathbf{b}), \end{aligned}$$

where  $\psi$  are the parameters of the fine-tuned PLMs,  $\oplus$  is the concatenation operation,  $p_\psi(\mathbf{Y} | \mathbf{S})$  is the label distribution predicted by the fine-tuned PLMs,  $p_\theta(\mathbf{Y} | \mathbf{Z}, \mathbf{T})$  is the label distribution predicted by MGM, which is based on Equation (7),  $\mathbf{W}$  and  $\mathbf{b}$  are the parameters of the linear classifier. More details are given in Figure 3 and Appendix E.

## 5 Experiments

### 5.1 Research Questions

We explore the following research questions:

- (RQ1) Can MGM tackle disconnected components and label sparsity in media graphs for factuality and political bias detection tasks?
- (RQ2) How do the number of global similar nodes  $K$  and the trade-off hyper-parameter  $\eta$  affect the performance of MGM?
- (RQ3) How does the memory module affect the performance of MGM?
- (RQ4) How does MGM elevate the performance of PLMs when faced with the challenge of missing text in *Wikipedia* or *Articles*?

Model	Fact-2020			Bias-2020		
	Macro-F1	Accuracy	Average Recall	Macro-F1	Accuracy	Average Recall
Majority class	22.93 $\pm$ 0.00	52.43 $\pm$ 0.00	33.33 $\pm$ 0.00	19.18 $\pm$ 0.00	40.39 $\pm$ 0.00	33.33 $\pm$ 0.00
GCN	25.55 $\pm$ 0.94	52.55 $\pm$ 0.28	34.74 $\pm$ 0.49	38.58 $\pm$ 5.13	42.90 $\pm$ 4.81	41.48 $\pm$ 5.11
+ MGM	<b>43.05 <math>\pm</math> 2.03</b>	<b>53.37 <math>\pm</math> 1.00</b>	<b>43.42 <math>\pm</math> 1.53</b>	<b>42.77 <math>\pm</math> 1.09</b>	<b>45.23 <math>\pm</math> 1.70</b>	<b>43.80 <math>\pm</math> 3.19</b>
GAT	33.75 $\pm$ 3.12	54.18 $\pm$ 0.77	39.26 $\pm$ 2.12	41.22 $\pm$ 1.79	50.34 $\pm$ 0.78	48.06 $\pm$ 1.05
+ MGM	<b>43.63 <math>\pm</math> 2.80</b>	<b>55.11 <math>\pm</math> 1.44</b>	<b>43.54 <math>\pm</math> 2.71</b>	<b>50.41 <math>\pm</math> 2.86</b>	<b>54.06 <math>\pm</math> 1.98</b>	<b>51.96 <math>\pm</math> 0.79</b>
GraphSAGE	42.68 $\pm$ 2.55	58.02 $\pm$ 1.18	45.70 $\pm$ 1.25	39.35 $\pm$ 1.07	50.00 $\pm$ 1.32	49.09 $\pm$ 1.06
+ MGM	<b>46.67 <math>\pm</math> 1.58</b>	<b>59.00 <math>\pm</math> 1.00</b>	<b>47.40 <math>\pm</math> 1.67</b>	<b>46.77 <math>\pm</math> 1.82</b>	<b>51.04 <math>\pm</math> 0.67</b>	<b>50.18 <math>\pm</math> 0.92</b>
SGC	22.73 $\pm$ 0.07	51.39 $\pm$ 0.28	33.10 $\pm$ 0.18	35.37 $\pm$ 0.60	45.34 $\pm$ 0.97	45.80 $\pm$ 0.76
+ MGM	<b>41.28 <math>\pm</math> 1.42</b>	<b>53.95 <math>\pm</math> 0.77</b>	<b>41.32 <math>\pm</math> 1.22</b>	<b>39.11 <math>\pm</math> 0.51</b>	<b>46.74 <math>\pm</math> 0.78</b>	<b>47.10 <math>\pm</math> 0.74</b>
DNA	22.75 $\pm$ 0.03	<b>51.74 <math>\pm</math> 0.00</b>	33.33 $\pm$ 0.00	24.27 $\pm$ 3.02	40.69 $\pm$ 0.73	35.03 $\pm$ 1.02
+ MGM	<b>34.04 <math>\pm</math> 1.60</b>	50.81 $\pm$ 1.30	<b>36.56 <math>\pm</math> 1.98</b>	<b>33.22 <math>\pm</math> 1.13</b>	<b>42.55 <math>\pm</math> 2.50</b>	<b>38.59 <math>\pm</math> 1.81</b>
FiLM	43.32 $\pm$ 2.25	57.09 $\pm$ 0.77	44.46 $\pm$ 1.40	39.33 $\pm$ 2.76	47.55 $\pm$ 1.12	47.85 $\pm$ 1.07
+ MGM	<b>49.68 <math>\pm</math> 1.62</b>	<b>57.90 <math>\pm</math> 2.39</b>	<b>49.94 <math>\pm</math> 1.68</b>	<b>45.33 <math>\pm</math> 2.76</b>	<b>48.25 <math>\pm</math> 2.65</b>	<b>48.61 <math>\pm</math> 2.84</b>
FAGCN	24.77 $\pm$ 7.52	47.04 $\pm$ 3.71	36.12 $\pm$ 5.30	19.69 $\pm$ 0.65	39.88 $\pm$ 0.28	33.71 $\pm$ 0.31
+ MGM	<b>48.77 <math>\pm</math> 0.00</b>	<b>53.14 <math>\pm</math> 1.66</b>	<b>49.19 <math>\pm</math> 0.00</b>	<b>45.02 <math>\pm</math> 3.00</b>	<b>45.69 <math>\pm</math> 2.88</b>	<b>45.07 <math>\pm</math> 3.00</b>
GATv2	51.42 $\pm$ 2.32	61.13 $\pm$ 1.04	55.36 $\pm$ 1.74	48.48 $\pm$ 1.68	55.11 $\pm$ 1.85	53.07 $\pm$ 1.75
+ MGM	<b>54.50 <math>\pm</math> 2.55</b>	<b>62.72 <math>\pm</math> 1.01</b>	<b>57.36 <math>\pm</math> 1.06</b>	<b>52.41 <math>\pm</math> 2.85</b>	<b>55.46 <math>\pm</math> 2.45</b>	<b>54.00 <math>\pm</math> 2.61</b>

Table 1: Performance of GNN baselines and their MGM enhanced versions for the factuality and political bias tasks on the ACL-2020 dataset, with the majority class baseline and SVM included as naïve and non-graphical methods. The higher performance is highlighted in **Bold**.

Model	Fact-2020		Bias-2020	
	Macro-F1 $\dagger/\$$	Average Recall $\dagger/\$$	Macro-F1 $\dagger/\$$	Average Recall $\dagger/\$$
GCN	42.04 $\pm$ 1.91 / <b>43.05 <math>\pm</math> 2.04</b>	41.97 $\pm$ 1.77 / <b>43.42 <math>\pm</math> 1.53</b>	<b>42.77 <math>\pm</math> 1.10</b> / 42.37 $\pm$ 2.09	43.72 $\pm$ 3.15 / <b>43.80 <math>\pm</math> 3.19</b>
GAT	40.63 $\pm$ 3.28 / <b>43.63 <math>\pm</math> 2.81</b>	40.23 $\pm$ 2.88 / <b>43.54 <math>\pm</math> 2.71</b>	<b>50.41 <math>\pm</math> 2.86</b> / 47.12 $\pm$ 1.69	<b>51.96 <math>\pm</math> 0.79</b> / 49.39 $\pm$ 2.02
GraphSAGE	45.11 $\pm$ 1.44 / <b>46.68 <math>\pm</math> 1.59</b>	45.69 $\pm$ 1.42 / <b>47.40 <math>\pm</math> 1.67</b>	44.96 $\pm$ 1.72 / <b>46.78 <math>\pm</math> 1.83</b>	49.78 $\pm$ 0.66 / <b>51.04 <math>\pm</math> 0.67</b>
SGC	<b>41.29 <math>\pm</math> 1.42</b> / 39.65 $\pm$ 2.00	<b>41.32 <math>\pm</math> 1.22</b> / 40.04 $\pm$ 1.50	38.32 $\pm$ 1.41 / <b>39.12 <math>\pm</math> 0.51</b>	46.74 $\pm$ 0.56 / <b>47.10 <math>\pm</math> 0.74</b>
DNA	33.48 $\pm$ 3.69 / <b>34.05 <math>\pm</math> 1.60</b>	35.99 $\pm$ 2.24 / <b>36.56 <math>\pm</math> 1.98</b>	<b>33.22 <math>\pm</math> 1.13</b> / 32.25 $\pm$ 4.14	<b>38.59 <math>\pm</math> 1.81</b> / 36.63 $\pm$ 4.53
FiLM	45.12 $\pm$ 3.38 / <b>49.68 <math>\pm</math> 1.62</b>	45.58 $\pm$ 2.78 / <b>49.94 <math>\pm</math> 1.68</b>	43.98 $\pm$ 2.57 / <b>45.33 <math>\pm</math> 2.76</b>	47.86 $\pm$ 1.46 / <b>48.61 <math>\pm</math> 2.84</b>
FAGCN	<b>48.77 <math>\pm</math> 0.00</b> / 46.88 $\pm$ 2.88	<b>49.19 <math>\pm</math> 0.00</b> / 48.05 $\pm$ 2.65	<b>45.02 <math>\pm</math> 3.00</b> / 44.36 $\pm$ 1.24	<b>45.07 <math>\pm</math> 3.00</b> / 44.47 $\pm$ 1.36
GATv2	54.13 $\pm$ 2.93 / <b>54.50 <math>\pm</math> 2.55</b>	56.82 $\pm$ 1.94 / <b>57.36 <math>\pm</math> 1.06</b>	<b>52.41 <math>\pm</math> 2.85</b> / 50.44 $\pm$ 0.95	<b>54.00 <math>\pm</math> 2.61</b> / 52.02 $\pm$ 1.29

Table 2: Summary of the MGM results detailing the performance variation the use of between using full memory ( $\dagger$ ) and a reduced (90%) memory allocation ( $\$$ ) for each GNN.

## 5.2 Dataset

The dataset for media factuality and political bias introduced by Baly et al. (2020c) comprises 859 media sources<sup>3</sup>, their domain names, and corresponding gold labels. These labels are sourced from MBFC, a platform supported by independent journalists. Factuality is given on a three-point scale: high, mixed, and low. Political bias is also on a three-point scale: left, center, right. Panayotov et al. (2022) used Alexa Rank<sup>4</sup> to create a graph based on audience overlap, using the 859 media as seed nodes. Media sources that shared the same audience, as determined by Alexa Rank, were connected with an edge. Alexa Rank returned a maximum of five similar media sources for each medium, which could be part of the initial seed nodes or newly identified media. In the resulting

graph, the nodes represent the media sources, and the edges represent the percentage of audience overlap between two media. We used these publicly available graph data (the only one of its kind) to train GNNs for the factuality and political bias of the media. More details are given in Appendix A.

## 5.3 Baselines

For evaluation, we consider two categories of baselines, including GNN-based and PLM-based models. For GNN models, we select eight well-known models, including GCN (Kipf and Welling, 2016), GraphSAGE (Hamilton et al., 2017), GAT (Veličković et al., 2018), SGC (Wu et al., 2019), DNA (Fey, 2019), FiLM (Brockschmidt, 2020), FAGCN (Bo et al., 2021) and GATv2 (Brody et al., 2022). More detail about these GNN baselines are provided in Appendix B. For PLMs, we use four popular encoder models, including BERT, RoBERTa, DistillBERT, and DeBERTaV3. We

<sup>3</sup><https://github.com/ramybaly/News-Media-Reliability>

<sup>4</sup><http://www.alexa.com/siteinfo>

Model	60% labels	80% labels	100% labels
GAT	37.90 $\pm$ 0.41	39.22 $\pm$ 0.71	33.75 $\pm$ 3.12
<b>+MGM</b>	<b>40.80 <math>\pm</math> 3.83</b>	<b>41.99 <math>\pm</math> 1.44</b>	<b>43.63 <math>\pm</math> 2.80</b>
FiLM	39.89 $\pm$ 1.69	38.59 $\pm$ 4.30	43.32 $\pm$ 2.25
<b>+MGM</b>	<b>42.93 <math>\pm</math> 4.00</b>	<b>38.62 <math>\pm</math> 2.73</b>	<b>49.68 <math>\pm</math> 1.62</b>
FAGCN	24.32 $\pm$ 3.18	22.73 $\pm$ 0.00	24.77 $\pm$ 7.52
<b>+MGM</b>	<b>34.10 <math>\pm</math> 7.36</b>	<b>39.89 <math>\pm</math> 4.04</b>	<b>48.77 <math>\pm</math> 0.00</b>
GATv2	40.54 $\pm$ 1.47	42.14 $\pm$ 2.76	51.42 $\pm$ 2.32
<b>+MGM</b>	<b>42.98 <math>\pm</math> 2.51</b>	<b>44.35 <math>\pm</math> 2.14</b>	<b>54.50 <math>\pm</math> 2.55</b>

Table 3: The impact of different proportions of training labeled data on the performance (Macro-F1) of MGM for the Fact-2020 task.

compare our results to state-of-the-art results for the factuality and political bias of the news media (Panayotov et al., 2022; Mehta et al., 2022).

## 6 Discussion

### 6.1 Overall Performance

To answer **RQ1**, we conduct factuality and political bias classification experiments in a semi-supervised setting. The experimental results reported in Table 1 demonstrate that MGM can improve the performance of existing GNNs in almost all cases. For example, when applied on the Fact-2020 dataset, MGM improves the Macro-F1 performance of GCN, GAT, SGC, and DNA by 17.5%, 9.8%, 18.5%, and 11.4%, respectively. Similarly, for Bias-2020, we can observe that GNNs equipped with MGM consistently outperform the corresponding base models in all evaluation measures.

We conducted a series of experiments using different proportions of training labels to assess the performance of MGM as shown in Table 3. The results indicate a clear trend: as we increase the percentage of training labels, the model’s performance improves significantly compared to the baseline. Due to limited data, using a smaller percentage of training labels results in modest improvements over the baseline, constraining the model’s ability to generalize well to unseen data. MGM effectively addresses **RQ1** by leveraging global similar nodes in media graphs with disconnected components and label sparsity for factuality and political bias detection. Our evaluation extends to Fact-2018 and depicts MGM’s stable performance across different datasets presented in Appendix D.

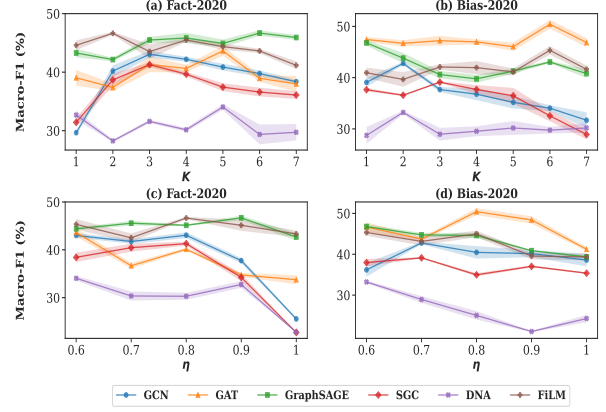


Figure 2: MGM performance across all GNNs for both tasks, evaluated for different values of  $K$  (global similar nodes) and  $\eta$  (trade-off hyper-parameter).

### 6.2 Impact of the Number of Global Similar Nodes

Next, we turn to **RQ2** to understand the impact of the number of global similar nodes  $K$ . Specifically, we investigate the performance of MGM with different values of  $K$ . As shown in Figures 2(a) and 2(b), leveraging a few global similar nodes can improve the performance of the base GNNs. For example, both GCN and SGC exhibit similar patterns, peaking in performance at  $K=3$  on the factuality task. The performance of GNNs enhanced with MGM decreases when  $K$  exceeds a certain threshold. This is attributed to the introduction of noise by incorporating excessive information from numerous global similar nodes.

### 6.3 Impact of the Trade-off Hyper-Parameter

Recall that in Section 4.1, we introduced a hyper-parameter  $\eta$  as a hyper-parameter that influences the predicted label distribution. When  $\eta = 1$ , MGM only relies on local node representations for prediction, degrading to a vanilla GNN. In contrast, when  $\eta < 1$ , our model incorporates information from both local neighbors and global similar nodes to predict the node labels. To further investigate the impact of the trade-off hyper-parameter  $\eta$ , we analyze the sensitivity of MGM to its value. The experimental results are shown in Figures 2(c) and 2(d). We find that compared to  $\eta = 1$ , MGM yields improved performance when  $\eta < 1$  in most cases. For example, GCN achieves its best performance when integrated with MGM using an  $\eta$  value of 0.8. Hence, the effectiveness of incorporating information from global similar nodes highlighted in the results validates the **RQ2**.

Model	Fact-2020						Bias-2020					
	Articles			Wikipedia			Articles			Wikipedia		
	Macro-F1	Accuracy	Avg Recall	Macro-F1	Accuracy	Avg Recall	Macro-F1	Accuracy	Avg Recall	Macro-F1	Accuracy	Avg Recall
STAGE 1	BERT <sub>Base</sub>	<b>38.27</b>	<b>63.37</b>	<b>39.65</b>	34.64	59.30	37.98	<b>65.38</b>	<b>68.02</b>	<b>64.01</b>	58.70	62.79
	RoBERTa <sub>Base</sub>	33.55	62.79	37.65	25.29	59.30	33.36	63.34	65.70	62.41	58.73	62.78
	DistilBERT <sub>Base</sub>	35.27	62.21	37.65	25.27	61.01	33.30	65.04	67.44	63.91	58.71	62.85
	DeBERTaV3 <sub>Base</sub>	25.28	61.02	33.35	<b>40.81</b>	<b>61.05</b>	<b>40.75</b>	58.72	62.80	58.69	<b>59.65</b>	<b>63.37</b>
STAGE 2	BERT <sub>MGM<sub>GATv2</sub></sub>	<b>76.18</b>	<b>81.98</b>	<b>71.86</b>	73.69	81.40	70.92	83.74	84.30	83.88	<b>82.25</b>	<b>82.56</b>
	RoBERTa <sub>MGM<sub>GATv2</sub></sub>	69.89	80.23	66.85	72.73	79.65	70.52	85.51	86.05	85.38	81.32	81.98
	DistilBERT <sub>MGM<sub>GATv2</sub></sub>	74.55	81.40	71.48	73.03	80.23	70.84	87.20	87.79	86.90	80.68	81.40
	DeBERTaV3 <sub>MGM<sub>GATv2</sub></sub>	64.87	77.91	62.97	<b>74.56</b>	<b>81.98</b>	<b>73.20</b>	<b>87.71</b>	<b>88.37</b>	<b>87.70</b>	80.26	80.81

Table 4: **Stage 1:** Performance of logistic regression (meta-learner) on PLM probabilities with missing media attributed as probabilities (0.0, 0.0, 0.0). **Stage 2:** Performance of the logistic regression (meta-learner) on PLMs probabilities + MGM<sub>GATv2</sub> probabilities for missing media for factuality and political bias of the ACL-2020 dataset.

Model	Fact-2020			Bias-2020		
	Macro-F1	Accuracy	Avg Recall	Macro-F1	Accuracy	Avg Recall
Node classification (NC) (Mehta et al., 2022)	68.90	63.72	-	-	-	-
InfOP Best Model (Mehta et al., 2022)	72.55	66.89	-	-	-	-
GRENNER (Panayotov et al., 2022)	69.61	74.27	-	91.93	92.08	-
STAGE 3	DeBERTaV3 <sub>MGM<sub>GATv2</sub></sub> + BERT <sub>MGM<sub>GATv2</sub></sub>	78.43	83.04	75.03	92.64	92.98
STAGE 4	DeBERTaV3 <sub>MGM<sub>GATv2</sub></sub> + BERT <sub>MGM<sub>GATv2</sub></sub> + MGM <sub>FILM</sub>	<b>79.72 ± 0.00</b>	<b>84.21 ± 0.00</b>	<b>76.54 ± 0.00</b>	93.04 ± 0.26	<b>93.45 ± 0.23</b>
	DeBERTaV3 <sub>MGM<sub>GATv2</sub></sub> + BERT <sub>MGM<sub>GATv2</sub></sub> + MGM <sub>FAGCN</sub>	75.69 ± 3.49	81.29 ± 3.09	72.24 ± 3.35	<b>93.08 ± 0.24</b>	<b>93.45 ± 0.23</b>
	DeBERTaV3 <sub>MGM<sub>GATv2</sub></sub> + BERT <sub>MGM<sub>GATv2</sub></sub> + MGM <sub>GATv2</sub>	77.96 ± 0.30	82.69 ± 0.29	74.84 ± 0.16	92.71 ± 0.46	93.10 ± 0.44

Table 5: Previous studies (Mehta et al., 2022; Panayotov et al., 2022) and our best results. **Stage 3:** We concatenate the probabilities of the best PLMs from *Wikipedia* and *Articles* and use logistic regression to make predictions. **Stage 4:** We use probabilities from the Stage 3 model and concatenate with the probabilities of three GNNs (MGM<sub>FILM</sub>, MGM<sub>FAGCN</sub> and MGM<sub>GATv2</sub>).

## 6.4 Effectiveness of the Memory Module

Recall that in Section 4.3, MGM leveraged a Dirichlet prior to select a small set of candidate nodes and stored their node embeddings in the sampled memory ( $M_S$ ) for prediction. To compare the effectiveness of the sampled memory module to the full memory module ( $M_G$ ), which stores all the training node embeddings, we conducted a performance comparison between the two memory modules. The experimental results are given in Table 2, and they answer **RQ3** that MGM using sampled memory achieves comparable performance to MGM when using full memory. For example, for the GAT model, the performance is higher when using sampled memory compared to using full memory. This suggests that the sampled memory effectively captures sufficient information, allowing MGM to maintain its performance even with limited memory. Experimental results on the Fact-2018 dataset are reported in the Appendix D.

## 6.5 Impact of Integrating MGM with PLMs

To answer **RQ4**, we integrate the MGM probabilities with those from deep learning models based on textual features, and we observe that this substantially enhances the performance. Initially, with zero probabilities for missing textual features, we achieved accuracies of 68.02% for political bias

in *Articles*, 63.37% for *Wikipedia*, 63.37% for factuality in *Articles*, and 61.05% for *Wikipedia*, respectively (see Table 4). Replacing the zero probabilities with the best MGM<sub>GATv2</sub> improved the performance by up to 30%. Further concatenating the best model probabilities in stage three led to additional gains, and in stage four, our models outperformed previous state-of-the-art results in political bias and factuality (Panayotov et al., 2022; Mehta et al., 2022) (see Table 5).

## 7 Conclusion and Future Work

We introduced MGM, an innovative EM framework that substantially improves the performance of GNNs on media graphs by leveraging global similar nodes. The external memory module of MGM efficiently stores and retrieves node representations, addressing the challenge of test-time inefficiency by selecting global similar nodes from a smaller candidate set based on a sparse node selection distribution. Our experiments demonstrate that the integration of MGM features with PLMs consistently improves over existing baselines and establishes new state-of-the-art results.

In future work, we plan to explore multi-graph fusion, multi-task learning, and ordinal classification for diverse graph structures in media profiling.



## Limitations

The graph dataset, originating from ACL-2020 media nodes, was constructed using the Alexa Rank siteinfo tool, which is currently unavailable. Although the graph aids the task by capturing the inherent and hidden relationships between media, building such graphs is complex and resource-intensive. The research largely relies on Western definitions of political bias (left/center/right), which may not accurately capture the nuanced ideological biases present in news outlets from other cultural or political contexts. Moreover, the available graph is limited to the 2020 dataset. We are actively working on constructing graphs for the latest benchmarks, which include a larger number of media sources and updated MBFC rankings. We faced limitations in collecting *Articles* and *Wikipedia* texts from media sources from the ACL-2020 dataset due to the inaccessibility of their websites.

## Ethical Statement

Optimizing model architectures to enhance energy efficiency in training and inference operations is crucial for reducing environmental impact. Instead of relying on extensive computational resources to train complex models—which significantly increase carbon emissions, we propose improving model performance with less computational power. The *Articles* from the media pages were compiled in strict compliance with legal and ethical standards. We carefully reviewed the terms of use for all websites to ensure that our data collection processes adhered to them. Our compilation focused solely on publicly available data, avoiding paywalls and subscription models. Transparent data collection methods were designed to minimize the impact on source websites, including limiting the access frequency to prevent resource strain.

## References

- Sadia Afroz, Michael Brennan, and Rachel Greenstadt. 2012. Detecting hoaxes, frauds, and deception in writing style online. In *2012 IEEE Symposium on Security and Privacy*, pages 461–475. IEEE.
- Jisun An, Meeyoung Cha, Krishna Gummadi, Jon Crowcroft, and Daniele Quercia. 2012. Visualizing media bias through Twitter. In *AAAI ICWSM*, volume 6.
- Jisun An, Meeyoung Cha, P. Krishna Gummadi, and Jon Crowcroft. 2011. Media landscape in Twitter: A

world of new conventions and political diversity. In *AAAI ICWSM*.

- Dilshod Azizov, S Liang, and P Nakov. 2023. Frank at checkthat! 2023: Detecting the political bias of news articles and news media. *Working Notes of CLEF*.

- Ramy Baly, Giovanni Da San Martino, James Glass, and Preslav Nakov. 2020a. [We can detect your bias: Predicting the political ideology of news articles](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4982–4991, Online. Association for Computational Linguistics.

- Ramy Baly, Georgi Karadzhov, Dimitar Alexandrov, James Glass, and Preslav Nakov. 2018. [Predicting factuality of reporting and bias of news media sources](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3528–3539, Brussels, Belgium. Association for Computational Linguistics.

- Ramy Baly, Georgi Karadzhov, Jisun An, Haewoon Kwak, Yoan Dinkov, Ahmed Ali, James Glass, and Preslav Nakov. 2020b. What was written vs. who read it: news media profiling using text analysis and social media context. *arXiv preprint arXiv:2005.04518*.

- Ramy Baly, Georgi Karadzhov, Jisun An, Haewoon Kwak, Yoan Dinkov, Ahmed Ali, James Glass, and Preslav Nakov. 2020c. [What was written vs. who read it: News media profiling using text analysis and social media context](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3364–3374, Online. Association for Computational Linguistics.

- Alberto Barrón-Cedeño, Firoj Alam, Tommaso Caselli, Giovanni Da San Martino, Tamer Elsayed, Andrea Galassi, Fatima Haouari, Federico Ruggeri, Julia Maria Struß, Rabindra Nath Nandi, et al. 2023a. The clef-2023 checkthat! lab: Checkworthiness, subjectivity, political bias, factuality, and authority. In *European Conference on Information Retrieval*, pages 506–517. Springer.

- Alberto Barrón-Cedeño, Firoj Alam, Andrea Galassi, Giovanni Da San Martino, Preslav Nakov, Tamer Elsayed, Dilshod Azizov, Tommaso Caselli, Gullal S Cheema, Fatima Haouari, et al. 2023b. Overview of the clef-2023 checkthat! lab on checkworthiness, subjectivity, political bias, factuality, and authority of news articles and their source. In *International Conference of the Cross-Language Evaluation Forum for European Languages*, pages 251–275. Springer.

- Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. 2018. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*.

694	Deyu Bo, Xiao Wang, Chuan Shi, and Huawei Shen.	Giovanni Da San Martino, Firoj Alam, Maram Hasanain,	746
695	2021. Beyond low-frequency information in graph	Rabindra Nath Nandi, Dilshod Azizov, and Preslav	747
696	convolutional networks. In <i>Proceedings of the AAAI</i>	Nakov. 2023. Overview of the CLEF-2023 Check-	748
697	<i>Conference on Artificial Intelligence</i> , volume 35,	That! lab task 3 on political bias of news arti-	749
698	pages 3950–3957.	cles and news media. In <i>Working Notes of CLEF</i>	750
699	Cristian Bodnar, Fabrizio Frasca, Yuguang Wang, Nina	2023–Conference and Labs of the Evaluation Forum,	751
700	Otter, Guido F Montufar, Pietro Lio, and Michael	CLEF '2023, Thessaloniki, Greece.	752
701	Bronstein. 2021. Weisfeiler and lehman go topolog-	Enyan Dai, Charu Aggarwal, and Suhang Wang. 2021.	753
702	ical: Message passing simplicial networks. In <i>Inter-</i>	Nrgnn: Learning a label noise resistant graph neural	754
703	<i>national Conference on Machine Learning</i> , pages	network on sparsely and noisily labeled graphs. In	755
704	1026–1037.	<i>Proceedings of the 27th ACM SIGKDD Conference</i>	756
705	Aleksandar Bojchevski and Stephan Günnemann. 2018.	<i>on Knowledge Discovery &amp; Data Mining</i> , pages 227–	757
706	Deep gaussian embedding of graphs: Unsupervised	236.	758
707	inductive learning via ranking. In <i>International Con-</i>	Enyan Dai, Wei Jin, Hui Liu, and Suhang Wang. 2022.	759
708	<i>ference on Learning Representations</i> .	Towards robust graph neural networks for noisy	760
709	Marc Brockschmidt. 2020. Gnn-film: Graph neural	graphs with sparse labels. In <i>Proceedings of the Fif-</i>	761
710	networks with feature-wise linear modulation. In <i>Inter-</i>	<i>teenth ACM International Conference on Web Search</i>	762
711	<i>national Conference on Machine Learning</i> , pages	<i>and Data Mining</i> , pages 181–191.	763
712	1144–1152. PMLR.	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	764
713	Shaked Brody, Uri Alon, and Eran Yahav. 2021. How	Kristina Toutanova. 2018. Bert: Pre-training of deep	765
714	attentive are graph attention networks? In <i>Internat-</i>	bidirectional transformers for language understand-	766
715	<i>ional Conference on Learning Representations</i> .	ing. <i>arXiv preprint arXiv:1810.04805</i> .	767
716	Shaked Brody, Uri Alon, and Eran Yahav. 2022. <a href="#">How</a>	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	768
717	<a href="#">attentive are graph attention networks?</a>	Kristina Toutanova. 2019. <a href="#">BERT: Pre-training of</a>	769
718	Ceren Budak, Sharad Goel, and Justin M Rao. 2016.	<a href="#">deep bidirectional transformers for language under-</a>	770
719	Fair and balanced? quantifying media bias through	<a href="#">standing</a> . In <i>Proceedings of the 2019 Conference of</i>	771
720	crowdsourced content analysis. <i>Public Opinion</i>	<i>the North American Chapter of the Association for</i>	772
721	<i>Quarterly</i> , 80(S1):250–271.	<i>Computational Linguistics: Human Language Tech-</i>	773
722	Sonia Castelo, Thais Almeida, Anas Elghafari, Aécio	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	774
723	Santos, Kien Pham, Eduardo Nakamura, and Juliana	4171–4186, Minneapolis, Minnesota. Association for	775
724	Freire. 2019. A topic-agnostic approach for identi-	Computational Linguistics.	776
725	fying fake news pages. In <i>Companion proceedings</i>	James Fairbanks, Natalie Fitch, Nathan Knauf, and Er-	777
726	<i>of the 2019 World Wide Web conference</i> , pages 975–	ica Briscoe. 2018. Credibility assessment in the news:	778
727	980.	do we need to read. In <i>Proc. of the MIS2 Work-</i>	779
728	Wei-Fan Chen, Henning Wachsmuth, Khalid Al Khatib,	<i>shop held in conjunction with 11th Int'l Conf. on Web</i>	780
729	and Benno Stein. 2018. Learning to flip the bias of	<i>Search and Data Mining</i> , pages 799–800. ACM.	781
730	news headlines. In <i>Proceedings of the 11th Interna-</i>	Lisa Fan, Marshall White, Eva Sharma, Ruisi Su, Pra-	782
731	<i>tional conference on natural language generation</i> ,	fulla Kumar Choubey, Ruihong Huang, and Lu Wang.	783
732	pages 79–88.	2019. In plain sight: Media bias through the lens of	784
733	Djork-Arné Clevert, Thomas Unterthiner, and Sepp	factual reporting. <i>arXiv preprint arXiv:1909.02670</i> .	785
734	Hochreiter. 2016. Fast and accurate deep network	Matthias Fey. 2019. Just jump: Dynamic neighborhood	786
735	learning by exponential linear units (elus). In <i>Inter-</i>	aggregation in graph neural networks. <i>arXiv preprint</i>	787
736	<i>national Conference on Learning Representations</i> .	<i>arXiv:1904.04849</i> .	788
737	Nadia K Conroy, Victoria L Rubin, and Yimin Chen.	Matthias Fey, Jan E Lenssen, Frank Weichert, and Jure	789
738	2015. Automatic deception detection: Methods for	Leskovec. 2021. Gnnautoscale: Scalable and expres-	790
739	finding fake news. <i>Proceedings of the association for</i>	sive graph neural networks via historical embeddings.	791
740	<i>information science and technology</i> , 52(1):1–4.	In <i>International Conference on Machine Learning</i> ,	792
741	Ganqu Cui, Jie Zhou, Cheng Yang, and Zhiyuan	pages 3294–3304.	793
742	Liu. 2020. Adaptive graph encoder for attributed	Matthias Fey and Jan Eric Lenssen. 2019. Fast graph	794
743	graph embedding. In <i>Proceedings of the 26th ACM</i>	representation learning with pytorch geometric. In	795
744	<i>SIGKDD international conference on knowledge dis-</i>	<i>ICLR 2019 (RLGM Workshop)</i> .	796
745	<i>covery &amp; data mining</i> , pages 976–985.	Xiaobo Guo, Weicheng Ma, and Soroush Vosoughi.	797
		2022. Measuring media bias via masked language	798
		modeling. In <i>Proceedings of the International AAAI</i>	799
		<i>Conference on Web and Social Media</i> , volume 16,	800
		pages 1404–1408.	801

802	Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017.	Nikhil Mehta and Dan Goldwasser. 2023. An interactive	856
803	Inductive representation learning on large graphs. <i>Ad-</i>	framework for profiling news media sources. <i>arXiv</i>	857
804	<i>advances in Neural Information Processing Systems</i> ,	<i>preprint arXiv:2309.07384</i> .	858
805	30.		
806	Junxian He, Taylor Berg-Kirkpatrick, and Graham Neu-	Nikhil Mehta, María Leonor Pacheco, and Dan Gold-	859
807	big. 2020. Learning sparse prototypes for text gener-	wasser. 2022. Tackling fake news detection by con-	860
808	ation. <i>Advances in Neural Information Processing</i>	tinually improving social context representations us-	861
809	<i>Systems</i> , 33:14724–14735.	ing graph neural networks. In <i>Proceedings of the</i>	862
810	Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2021.	<i>60th Annual Meeting of the Association for Computa-</i>	863
811	Debertav3: Improving deberta using electra-style pre-	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages	864
812	training with gradient-disentangled embedding shar-	1363–1380.	865
813	ing. <i>arXiv preprint arXiv:2111.09543</i> .		
814	Austin Hounsel, Jordan Holland, Ben Kaiser, Kevin	Subhabrata Mukherjee and Gerhard Weikum. 2015.	866
815	Borgolte, Nick Feamster, and Jonathan Mayer. 2020.	<a href="#">Leveraging joint interactions for credibility analy-</a>	867
816	Identifying disinformation websites using infrastruc-	<a href="#">sis in news communities</a> . In <i>Proceedings of the 24th</i>	868
817	ture features. In <i>10th USENIX Workshop on Free and</i>	<i>ACM International Conference on Information and</i>	869
818	<i>Open Communications on the Internet (FOCI 20)</i> .	<i>Knowledge Management, CIKM 2015, Melbourne,</i>	870
819	Minyoung Huh, Andrew Liu, Andrew Owens, and	<i>VIC, Australia, October 19 - 23, 2015</i> , pages 353–	871
820	Alexei A Efros. 2018. Fighting fake news: Image	362. ACM.	872
821	splice detection via learned self-consistency. In <i>Pro-</i>		
822	<i>ceedings of the European conference on computer</i>	Vinod Nair and Geoffrey E Hinton. 2010. Rectified	873
823	<i>vision (ECCV)</i> , pages 101–117.	linear units improve restricted boltzmann machines.	874
824	Seokho Kang. 2021. K-nearest neighbor learning with	In <i>Proceedings of the 27th International Conference</i>	875
825	graph neural networks. <i>Mathematics</i> , 9(8):830.	<i>on International Conference on Machine Learning</i> ,	876
826	Johannes Kiesel, Maria Mestre, Rishabh Shukla, Em-	pages 807–814.	877
827	manuel Vincent, Payam Adineh, David Corney,		
828	Benno Stein, and Martin Potthast. 2019. <a href="#">SemEval-</a>	Preslav Nakov, Firoj Alam, Giovanni Da San Martino,	878
829	<a href="#">2019 task 4: Hyperpartisan news detection</a> . In	Maram Hasanain, Rabindra Nath Nandi, Dilshod Az-	879
830	<i>Proceedings of the 13th International Workshop on</i>	izov, and Panayot Panayotov. 2023a. Overview of the	880
831	<i>Semantic Evaluation</i> , pages 829–839, Minneapolis,	CLEF-2023 CheckThat! lab task 4 on factuality of	881
832	Minnesota, USA. Association for Computational Lin-	reporting of news media. In <i>Working Notes of CLEF</i>	882
833	guistics.	<i>2023–Conference and Labs of the Evaluation Forum</i> ,	883
834	Diederik P Kingma and Jimmy Ba. 2015. Adam: A	CLEF ’2023, Thessaloniki, Greece.	884
835	method for stochastic optimization. In <i>International</i>		
836	<i>Conference on Learning Representations</i> .	Preslav Nakov, Firoj Alam, Giovanni Da San Martino,	885
837	Thomas N Kipf and Max Welling. 2016. Semi-	Maram Hasanain, RN Nandi, D Azizov, and P Panay-	886
838	supervised classification with graph convolutional	otov. 2023b. Overview of the clef-2023 checkthat!	887
839	networks. <i>arXiv preprint arXiv:1609.02907</i> .	lab task 4 on factuality of reporting of news media.	888
840	Vivek Kulkarni, Junting Ye, Steve Skiena, and	<i>Working Notes of CLEF</i> .	889
841	William Yang Wang. 2018. <a href="#">Multi-view models for</a>	Preslav Nakov, Husrev Taha Sencar, Jisun An, and Hae-	890
842	<a href="#">political ideology detection of news articles</a> . In <i>Pro-</i>	woon Kwak. 2021. A survey on predicting the fac-	891
843	<i>ceedings of the 2018 Conference on Empirical Meth-</i>	tuality and the bias of news media. <i>arXiv preprint</i>	892
844	<i>ods in Natural Language Processing</i> , pages 3518–	<i>arXiv:2103.12506</i> .	893
845	3527, Brussels, Belgium. Association for Computa-		
846	tional Linguistics.	Panayot Panayotov, Utsav Shukla, Husrev Taha Sen-	894
847	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	car, Mohamed Nabeel, and Preslav Nakov. 2022.	895
848	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	GREENER: Graph neural networks for news media	896
849	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	profiling. In <i>Proceedings of the 2022 Conference on</i>	897
850	Roberta: A robustly optimized bert pretraining ap-	<i>Empirical Methods in Natural Language Processing</i> ,	898
851	proach. <i>arXiv preprint arXiv:1907.11692</i> .	EMNLP ’22, Abu Dhabi, UAE.	899
852	Guixiang Ma, Vy Vo, Theodore Willke, and Nesreen K	Verónica Pérez-Rosas, Bennett Kleinberg, Alexandra	900
853	Ahmed. 2022. Memory-augmented graph neural net-	Lefevre, and Rada Mihalcea. 2017. Automatic detec-	901
854	works: A neuroscience perspective. <i>arXiv preprint</i>	tion of fake news. <i>arXiv preprint arXiv:1708.07104</i> .	902
855	<i>arXiv:2209.10818</i> .		
		Kashyap Popat, Subhabrata Mukherjee, Jannik Strötgen,	903
		and Gerhard Weikum. 2017. Where the truth lies:	904
		Explaining the credibility of emerging claims on the	905
		Web and social media. In <i>WWW Companion</i> , pages	906
		1003–1012.	907
		Martin Potthast, Johannes Kiesel, Kevin Reinartz, Janek	908
		Bevendorff, and Benno Stein. 2018. <a href="#">A stylometric</a>	909
		<a href="#">inquiry into hyperpartisan and fake news</a> . In <i>Proceed-</i>	910
		<i>ings of the 56th Annual Meeting of the Association for</i>	911



*Computational Linguistics (Volume 1: Long Papers)*, pages 231–240, Melbourne, Australia. Association for Computational Linguistics.

Meng Qu, Yoshua Bengio, and Jian Tang. 2019. Gmnn: Graph markov neural networks. In *International Conference on Machine Learning*, pages 5241–5250.

Meng Qu, Huiyu Cai, and Jian Tang. 2021. Neural structured prediction for inductive node classification. In *International Conference on Learning Representations*.

Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.

Jayaram Sethuraman. 1994. A constructive definition of dirichlet priors. *Statistica sinica*, pages 639–650.

Timo Spinde, Manuel Plank, Jan-David Krieger, Terry Ruas, Bela Gipp, and Akiko Aizawa. 2022. Neural media bias detection using distant supervision with BABE–bias annotations by experts. *arXiv preprint arXiv:2209.14557*.

Peter Stefanov, Kareem Darwish, Atanas Atanasov, and Preslav Nakov. 2020. [Predicting the topical stance and political leaning of media using tweets](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 527–537, Online. Association for Computational Linguistics.

Ke Sun, Zhouchen Lin, and Zhanxing Zhu. 2020. Multi-stage self-supervised learning for graph convolutional networks on graphs with few labeled nodes. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 5892–5899.

Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph attention networks. In *International Conference on Learning Representations*.

Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. 2019. Simplifying graph convolutional networks. In *International Conference on Machine Learning*, pages 6861–6871.

Kai-Cheng Yang and Filippo Menczer. 2023. Large language models can rate news outlet credibility. *arXiv preprint arXiv:2304.00228*.

Mingqi Yang, Renjian Wang, Yanming Shen, Heng Qi, and Baocai Yin. 2022. Breaking the expression bottleneck of graph neural networks. *IEEE Transactions on Knowledge and Data Engineering*.

Ningyu Zhang, Xin Xie, Xiang Chen, Shumin Deng, Chuanqi Tan, Fei Huang, Xu Cheng, and Huajun Chen. 2022. Reasoning through memorization: Nearest neighbor knowledge graph embeddings. *arXiv preprint arXiv:2201.05575*.



## Appendix

### A GNN Data & Task Statistics

Table 6 describes the statistics of the graph data. Factuality is given on a three-point scale: high, mixed, and low. Political bias is also on a three-point scale: left, center, right. Panayotov et al. (2022) used the Alexa<sup>5</sup> (down temporarily) to create a graph based on audience overlap, using 859 media from ACL-2020 (Baly et al., 2020c) as seed nodes. Media sources that shared the same audience, as determined by Alexa, were connected with an edge, provided they met a specific score threshold. Alexa was set to return a maximum of five similar media sources for each medium; these could be part of the initial seed nodes or newly identified media. The primary graph constructed using ACL-2020 dataset media as seed nodes is designated as *level-0*. In this graph, the nodes represent the media sources that publish news or information, and the edges represent the audience overlap for a pair of nodes. The procedure was reiterated five times, leading to the formation of five distinct graph levels. With every subsequent iteration, the graph expanded, encompassing media sources previously identified by Alexa. This iterative expansion resulted in a progressive increase in both the number of nodes and edges at each level.

Upon analyzing the constructed graphs, we observed several disconnected components, each signifying a unique subnetwork of nodes. Naturally, as the graph levels increased, the number of these components decreased. This can be attributed to the fact that an increase in nodes offers more opportunities for components to merge. We opt for graph level 3 to train GNNs, as detailed in Table 6: it represents the most granular level publicly accessible with fewer disconnected components for both factual and bias tasks. The Alexa tool also generated features for each node in the graph, which we treat as node attributes while training the GNNs. These features include site rank, total sites linked in, bounce rate, and the daily time users spend on the site. These features are the numeric values that are described and normalized in the study (Panayotov et al., 2022). We refer to the GNN training tasks as *Fact-2020* and *Bias-2020* for the factuality and political bias tasks, respectively, since both tasks are derived using ACL-2020. As graph-based data become increasingly accessible,

<sup>5</sup><http://www.alexa.com/siteinfo>

Property	Specification
Nodes	67,350
Edges	200,481
Features	5
Discon. comp.	44
Avg. nodes / comp.	1,500
labeled Nodes	859 (1%)
Unlabeled Nodes	66,492 (99%)
Tasks	Fact-2020, Bias-2020
Factuality task dist.	high (162), mix (249), low (453)
Political Bias task dist.	left (243), center (272), right (349)
Training Split	687 (80% of 1%)
Test Split	172 (20% of 1%)

Table 6: Statistics about the level-3 graph constructed from ACL-2020 (Panayotov et al., 2022).

we focus exclusively on the graph and its inherent features, promoting an approach tailored to such structures. In contrast, (Panayotov et al., 2022) operates in a supervised setting and uses specialized textual features (e.g., Articles, Wikipedia, Twitter, and YouTube) that are not publicly available. The proposed MGM addresses the unique challenges of the media graph, offering solutions to the research questions described in the designated section 5.1.

Table 7 describes the statistics of the level-3 graph constructed from EMNLP -2018 (Panayotov et al., 2022) media in the same way explained in Section 5.2. The EMNLP-2018 dataset comprises 1,066 news outlets, rated on a 3-point scale for factuality (*high, mixed, low*) and a 7-point scale for political bias (*extreme-left, left, center-left, center, center-right, right, and extreme-right*) (Baly et al., 2018). A subsequent analysis (Baly et al., 2020c) identified that the labels *center-left* and *center-right* serve as vague intermediate categories, leading to their exclusion. Furthermore, to minimize subjectivity in the annotator decisions, the *extreme-left* and *extreme-right* categories were amalgamated into the *left* and *right* categories, respectively. This adjustment resulted in a simplified 3-point political bias scale (*left, center, right*) and reduced the dataset to 864 outlets as shown in Table 6, published in ACL-2020, which we consider as our main dataset in section 5.2.

### B Baselines

This section summarizes the baseline GNN models that we use as the backbone for our proposed MGM framework to enhance their learning capabilities in the presence of sparsity challenges.

**GCN (Kipf and Welling, 2016):** GCN simplifies the convolution operation to alleviate the problem

Property	Specification
Nodes	78429
Edges	232530
Features	5
Discon. comp.	88
Avg. nodes / comp.	911
labeled Nodes	1066 (1.35%)
Unlabeled Nodes	77363 (98.65%)
Tasks	Fact-2018
Factuality task dist.	high (265), mixed (268), low (542)
Training Split	852 (80% of 1.35%)
Test Split	214 (20% of 1.35%)

Table 7: Statistics about the level-3 graph constructed from EMNLP-2018.

Property	Specification
Tasks	Fact-2020, Bias-2020
Factuality task dist.	high (295), mix (119), low (58)
Political Bias task dist.	left (152), center (181), right (139)
Training Split	387
Test Split	85

Table 8: Statistics about *Articles* and *Wikipedia* collected from ACL-2020 (Panayotov et al., 2022) dataset.

Hyper-parameter	BERT	RoBERTa	DistilBERT	DeBERTaV3
Batch size	80	100	120	80
Max length	512	512	512	512
Epochs	3	4	5	5
Learning rate	2e-5	2e-5	2e-5	2e-5

Table 9: Experimental setup for PLMs

of overfitting and introduces a renormalization trick to solve the vanishing gradient problem. We set the number of hidden neurons to 16, and the number of layers to 2. ReLU (Nair and Hinton, 2010) is used as the activation function. We do not dropout between GNN layers.

**SGC** (Wu et al., 2019): SGC shows that the graph convolution in GNNs is actually Laplacian smoothing, which smooths the feature matrix so that nearby nodes have similar hidden representations. SGC removes the weight matrices and nonlinearities between layers. In our experiments, we set the number of hidden neurons to 256, the number of layers to 2, and the number of hops to 2. We do not dropout between GNN layers.

**GraphSAGE** (Hamilton et al., 2017): GraphSAGE learns embeddings of nodes in the network by sampling and aggregating features from nodes’ local neighborhoods. GraphSAGE has different variants based on different feature aggregators, and we adopt GraphSAGE with a mean-based aggregator as our baseline. In our experiments, we set the number of hidden neurons to 64, and the number of layers to 2. ELU (Clevert et al., 2016) is used as

the activation function. We do not dropout between GNN layers.

**GAT** (Veličković et al., 2018): GAT incorporates the attention mechanism into the propagation step, allowing each node to compute its hidden states by attending to its neighbors using self-attention and multi-head attention strategies. we set the number of hidden neurons to 128 per attention head and the number of layers to 3. The number of heads for each layer is set to 4, 4 and 6. ELU (Clevert et al., 2016) is used as the activation function. We do not dropout between GNN layers.

**DNA** (Fey, 2019): DNA leverages the jumping knowledge network to enhance the performance of GNNs. This approach enables selective and node-adaptive aggregation of neighboring embeddings, even when they have different localities within the graph. We set the number of hidden neurons to 128, the number of heads to 8, and the number of layers to 4. ReLU (Nair and Hinton, 2010) is used as the activation function. We set the dropout rate to 0.5 between GNN layers.

**FiLM** (Brockschmidt, 2020): FiLM learns embeddings of nodes in the network by training a linear message function that is conditioned on the features of the neighboring nodes. This allows FiLM to effectively capture and incorporate contextual information from the neighbors into the node embeddings. We set the number of hidden neurons to 320 and the number of layers to 4. We set the dropout rate to 0.1 between GNN layers.

**FAGCN** (Bo et al., 2021): FAGCN adopts a self-gating attention mechanism to learn the proportion of low-frequency and high-frequency signals. By adaptively modelling the frequency signals, FAGCN achieves enhanced expressive performance in capturing graph structure and features. We set the number of hidden neurons to 16, and the number of layers to 4. We set the dropout rate to 0.5 between GNN layers. **GATv2Conv** (Brody et al., 2021): GATv2 introduces a dynamic graph attention variant that reorders the internal operations, resulting in a significantly higher level of expressiveness compared to GAT. We set the number of hidden neurons to 64 per attention head and the number of layers to 3. ELU (Clevert et al., 2016) is used as the activation function. We do not dropout between GNN layers.

Model	Macro-F1 $\dagger/\S$	Average Recall $\dagger/\S$
GCN	<b>47.20</b> $\pm$ <b>1.54</b> / 46.52 $\pm$ 1.52	<b>48.13</b> $\pm$ <b>1.19</b> / 47.60 $\pm$ 1.16
GAT	<b>54.99</b> $\pm$ <b>4.14</b> / 53.65 $\pm$ 2.79	<b>57.15</b> $\pm$ <b>4.05</b> / 55.85 $\pm$ 2.27
GraphSage	46.54 $\pm$ 1.65 / <b>47.86</b> $\pm$ <b>1.38</b>	49.09 $\pm$ 0.87 / <b>50.91</b> $\pm$ <b>1.07</b>
SGC	44.60 $\pm$ 2.41 / <b>45.16</b> $\pm$ <b>2.29</b>	45.82 $\pm$ 0.73 / <b>46.03</b> $\pm$ <b>1.80</b>
DNA	<b>34.93</b> $\pm$ <b>3.95</b> / 33.88 $\pm$ 1.52	<b>36.71</b> $\pm$ <b>3.83</b> / 35.35 $\pm$ 1.68
FiLM	51.06 $\pm$ 2.24 / <b>51.47</b> $\pm$ <b>2.47</b>	51.35 $\pm$ 2.16 / <b>52.13</b> $\pm$ <b>2.01</b>

Table 10: Summary of the MGM results detailing performance variations between using full memory ( $\dagger$ ) and a reduced 90% memory allocation ( $\S$ ) for each GNN across Fact-2018 task. The best performance per base model is marked in boldface.

## C Experimental Settings

As mentioned in Section 4.2, MGM is trained using the variational EM, which iteratively maximizes the ELBO and the expectation of log-likelihood function through an E-step and an M-step. To optimize the model, we use the Adam optimizer (Kingma and Ba, 2015) with a learning rate of 0.001. The early stopping strategy is implemented with patience of 10 epochs. In each experiment, we train MGM for 50 iterations to obtain the results. In order to encourage a sparse node selection distribution, we set the Dirichlet hyper-parameter  $\alpha$  to 0.1. The hyper-parameter  $K$ , which determines the number of global similar nodes, is selected from the range  $[1, 7]$  through a tuning process. Its value is optimized to achieve the best performance in the validation set for the node classification task. Similarly, the trade-off hyper-parameter  $\eta$ , which strikes a balance between the utilization of local representations and the information from global similar nodes, is chosen from the range  $[0.6, 1]$  and is tuned to obtain the optimal performance in the validation set for the node classification task. The model is trained for 5 epochs using different random seeds and mean  $\pm$  standard deviation is reported. We use the GNN module implementations provided by PyTorch Geometric<sup>6</sup> (Fey and Lenssen, 2019).

**Evaluation Measures** We evaluate our frameworks using mean of three key measures: *Macro-F1*, *Accuracy*, and *Average Recall*. *Macro-F1* balances precision and recall for each class, ideal for imbalanced datasets. *Accuracy* measures overall correctness, while *Average Recall* highlights the model’s sensitivity to different classes. For GNNs experiment, we used an Nvidia 2080 Ti GPU, and for PLMs experiment, we used an NVIDIA A6000

<sup>6</sup>[https://github.com/pyg-team/pytorch\\_geometric/tree/master/examples](https://github.com/pyg-team/pytorch_geometric/tree/master/examples)

Model	Macro-F1	Average Recall
Majority class	22.47 $\pm$ 0.00	33.33 $\pm$ 0.00
SVM	41.78 $\pm$ 0.00	48.89 $\pm$ 0.00
GCN	48.63 $\pm$ 2.19	48.16 $\pm$ 2.49
+ MGM	<b>49.21</b> $\pm$ <b>1.54</b>	<b>51.13</b> $\pm$ <b>1.19</b>
GAT	46.63 $\pm$ 3.53	52.25 $\pm$ 4.20
+ MGM	<b>54.99</b> $\pm$ <b>4.14</b>	<b>57.15</b> $\pm$ <b>4.05</b>
GraphSAGE	41.77 $\pm$ 0.22	48.65 $\pm$ 0.22
+ MGM	<b>47.86</b> $\pm$ <b>1.38</b>	<b>50.91</b> $\pm$ <b>1.07</b>
SGC	41.06 $\pm$ 0.35	44.91 $\pm$ 0.43
+ MGM	<b>45.16</b> $\pm$ <b>2.29</b>	<b>46.03</b> $\pm$ <b>1.80</b>
DNA	28.24 $\pm$ 1.23	33.26 $\pm$ 1.03
+ MGM	<b>34.93</b> $\pm$ <b>3.95</b>	<b>36.71</b> $\pm$ <b>3.83</b>
FiLM	46.75 $\pm$ 0.79	50.92 $\pm$ 1.36
+ MGM	<b>51.47</b> $\pm$ <b>2.47</b>	<b>52.13</b> $\pm$ <b>2.01</b>

Table 11: Performance of GNN baselines and their MGM enhanced versions on the Fact task of EMNLP-2018, with the majority class baseline and SVM included as naive and non-graphical methods. The highest performance is highlighted in bold.

48GB GPU.

## D Results on Dataset EMNLP-2018

This section presents the results of the MGM model on the factuality task from the EMNLP-2018 dataset, supplementing our evaluation of MGM’s adaptability across various datasets and tasks. Table 11 shows that MGM is able to enhance GNNs for news media graphs. All the GNNs baselines are consistently improved with MGM. Given the reasons mentioned in Appendix A, we did not conduct experiments on the political bias task of EMNLP-2018.

The experimental results in Table 10 demonstrate that MGM using sampled memory achieves performance comparable to MGM with full memory. Using sampled memory, all baselines perform competitively compared to the use of full memory. This suggests that sampled memory effectively captures crucial information while mitigating computational and storage overhead compared to full memory.

## E Collecting Articles & Wikipedia

*Articles.* The article collection involves the following steps: (i) We obtained media sources from the ACL-2020 dataset. (ii) During the article link parsing, we parsed front-page article links from these

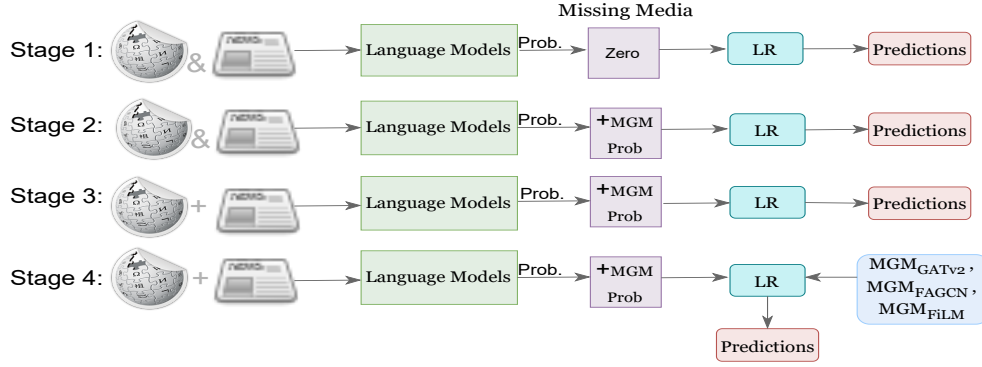


Figure 3: The pipeline of Integrating MGM with PLMs. **Stage 1:** We use logistic regression (meta-learner) to make predictions on probabilities obtained from PLMs on 472 media sources. For the remaining media sources, we assign  $[0.0, 0.0, 0.0]$  probabilities. **Stage 2:** We use the probabilities produced by PLMs and, for the missing ones, we integrate probabilities from the best GNN  $MGM_{GATv2}$ . Logistic regression is then used to make predictions. **Stage 3:** We concatenate the probabilities of the best PLMs on *Wikipedia* and *Articles* and use logistic regression to make predictions. **Stage 4:** We use the probabilities obtained from Stage 3, which involve concatenating these probabilities with those generated by three GNNs ( $MGM_{FILM}$ ,  $MGM_{FAGCN}$ , and  $MGM_{GATv2}$ ) across five different run seeds. Subsequently, logistic regression is employed to make predictions, and the scores are calculated using standard deviation.

media sources based on criteria of selecting only internal links with more than 65 characters and excluding menu button links. (iii) In the article collection stage, we use the selected article links to retrieve the titles and full text of the articles, using scripts and manual testing to ensure effective text extraction, with up to 30 news articles per media. (iv) Finally, the post-processing stage involved formatting the collected data in the JSON format. Moreover, we specifically targeted sections that focused on political, economic, and social issues sections.

*Wikipedia.* We started by searching for the media outlet name on the Internet to find the *Wikipedia* link. We ensured the link directed to a *Wikipedia* page specifically about the media outlet. We then retrieved the text from the *Wikipedia* page using its consistent HTML format. Finally, the post-processing stage involved formatting the collected data in the required JSON format.

In total, from 859 media sources, we have collected data from 472 media sources with *Articles* and *Wikipedia*. Table 8 provides detailed statistics. Moreover, Figure 3 provides our detailed pipeline for integrating MGM with PLMs.