

# Text Graph Neural Networks for Detecting AI-Generated Content

**Andric Valdez-Valenzuela**

Instituto de Investigaciones  
en Matemáticas Aplicadas  
y en Sistemas  
CDMX, Mexico.

**Helena Gómez-Adorno**

Instituto de Investigaciones  
en Matemáticas Aplicadas  
y en Sistemas  
CDMX, Mexico.

**Manuel Montes-y-Gómez**

Instituto Nacional  
de Astrofísica, Óptica  
y Electrónica  
Puebla, Mexico.

## Abstract

The widespread availability of Large Language Models (LLMs) such as GPT-4 and Llama-3, among others, has led to a surge in machine-generated content across various platforms, including social media, educational tools, and academic settings. While these models demonstrate remarkable capabilities in generating coherent text, their misuse raises significant concerns. For this reason, detecting machine-generated text has become a pressing need to mitigate these risks. This research proposed a novel classification method combining text-graph representations with Graph Neural Networks (GNNs) and different node feature initialization strategies to distinguish between human-written and machine-generated content. Experimental results demonstrate that the proposed approach outperforms traditional machine learning classifiers, highlighting the effectiveness of integrating structural and semantic relationships in text.

## 1 Introduction

Large language Models (LLMs) are now widely available and easily accessible, resulting in increased machine-generated content across various platforms, including Q&A forums, social media, educational resources, and academic contexts. Recent advancements in LLM technology, such as Llama-3 and GPT-4, have enabled these models to generate coherent responses to most user inquiries, making them increasingly attractive for replacing human labor in multiple fields. Moreover, this increased accessibility has led to concerns about misuse, including the creation of fake news, financial and legal issues, and education disruptions. Given the difficulty people have distinguishing between machine-generated and human-written text, there is a growing need for automated systems that can detect machine-generated content to mitigate and address the risks associated with its misuse (Nitu and Dascalu, 2024).

Viewing this problem as a classical text classification task, there are plenty of approaches to tackle it nowadays, from traditional methods such as training machine learning or deep learning models using Bag of Words or stylometric patterns as features to more advanced approaches based on the Transformer architecture, such as BERT (Devlin, 2018) or RoBERTa (Liu, 2019).

However, in recent years, a new area of research known as Graph Neural Networks (GNNs), or graph embeddings, has gained significant interest (Battaglia et al., 2018). These networks have proven to be highly effective in text classification tasks involving complex relational structures as they retain the global structure of a graph within their embeddings.

A text can be appropriately represented as a graph using the words/documents as nodes and the edge representing the significant relationship between the nodes. Also, it is possible to assign different attributes or weights to the graph's edges/nodes to add more significant information. Moreover, graph-based methods capture complex connections and dependencies that traditional methods might miss. Understanding these relationships between words/nodes is crucial for exploiting the text's best features. It could also help capture semantic and syntactic nuances in the text, distinguishing between human and machine-generated texts.

The contributions of our paper are summarized as follows:

- A classification method that combines Text-Graph representation and GNNs with different node feature initialization strategies to detect machine-generated text <sup>1</sup>.
- A detailed comparison of the performance of our method against baseline and state-of-the-art methods on English texts.

<sup>1</sup><https://github.com/andricValdez/GraphDeepLearning>



Additionally, the PMI weight measures the association strength between two words, highlighting how often these words co-occur more than would be expected by chance. Higher PMI values suggest a strong contextual relationship between the word pairs, even if their co-occurrence is infrequent. This metric helps reveal significant word associations that might not be immediately apparent from frequency alone.

### 3.2 Pipeline Architecture

Figure 2 shows the proposed pipeline architecture for identifying machine-generated text. As a first step, each document in the corpus is transformed into a Co-Occurrence graph, as described in Section 1. This transformation captures the relational structure of the texts, which is crucial for further processing.

After the graph has been built, we apply fine-tuning using pre-trained transformer models, such as BERT-Base-Uncased or RoBERTa-Base. These models are used to initialize node features in the graph, enhancing the semantic understanding of the text. Moreover, for the sake of comparison in the performance, we tested using different node features, such as the Word2Vec model (Mikolov, 2013) and random features.

The processed graph is then fed into a Graph Neural Network using a Graph Attention Network (GAT) layer (Veličković et al., 2017). This GAT layer uses attention mechanisms to focus on the most important nodes and edges, capturing relationships between words in a more subtle way. This attention-based learning enables the model to understand complex dependencies and associations within the text.

We implemented the GAT layer using the PyTorch Geometric library<sup>4</sup>, providing as inputs the node features, the co-occurrence graph (as a sparse matrix, in COO format) and the edge weights (freq or PMI metrics). Additionally, we set and test with different kinds of hyperparameter related to this GNN, such as the number of convolutions (message passing layers), head attentions, hidden channels, pooling layers (add, mean, max), and normalization layers (such as BatchNorm1d or Dropout).

Lastly, the graph document embedding is fed into a final classification model. This classifier, typically handled by a dense neural network (but it

<sup>4</sup>This library provides various methods for deep learning on graphs from a variety of published papers: <https://pytorch-geometric.readthedocs.io/en/latest/>

<i>Partition</i>	<i>Autextification 2023</i>		
	<b>human</b>	<b>machine</b>	<b>total</b>
<i>Train</i>	11,963	11,728	<b>23,691</b>
<i>Validation</i>	5,083	5,071	<b>10,154</b>
<i>Test</i>	10,642	11,190	<b>21,832</b>

Table 1: Summary stats for the Autextification 2023 English dataset used in the experiments.

could be any classification algorithm), determines whether the text was machine-generated based on the learned representations.

## 4 Experiments

This section shows all the experiment settings, datasets used, and the performance and results obtained for the proposed method.

### 4.1 Dataset

To evaluate the effectiveness of our proposed method, we utilized the Autextification2023 (Sarvazyan et al., 2023) dataset, a publicly available corpus specifically designed for machine-generated text detection. This dataset contains text in English and Spanish (In our case, we only used English texts). It comprises human and LLM-generated texts across five domains: tweets, reviews, how-to articles, news, and legal documents, representing a range of writing styles from formal to informal. Human texts were sourced from publicly available datasets like MultiEURLEX, XSUM, XL-SUM, MLSUM, Amazon Reviews, WikiLingua, and more. Machine-generated texts were produced using BLOOM and GPT-3 models, chosen for their multilingual capabilities and accessibility. Table 1 shows the Train, Validation, and Test sets, with a balanced distribution between human-generated and machine-generated text samples.

### 4.2 Results

Table 2 compares the performance of various models on the datasets for detecting machine-generated text using the Accuracy and F1-Score (macro) measures. The models evaluated include traditional machine learning classifiers, such as Support Vector Machine (using TF-IDF unigrams for vector representation), a fine-tuned BERT and RoBERTa model<sup>5</sup>, and GNNs with different node initialization strategies.

<sup>5</sup>Using a 16 batch size, five training epochs, 2e-5 as learning rate and 0.01 of weight decay.

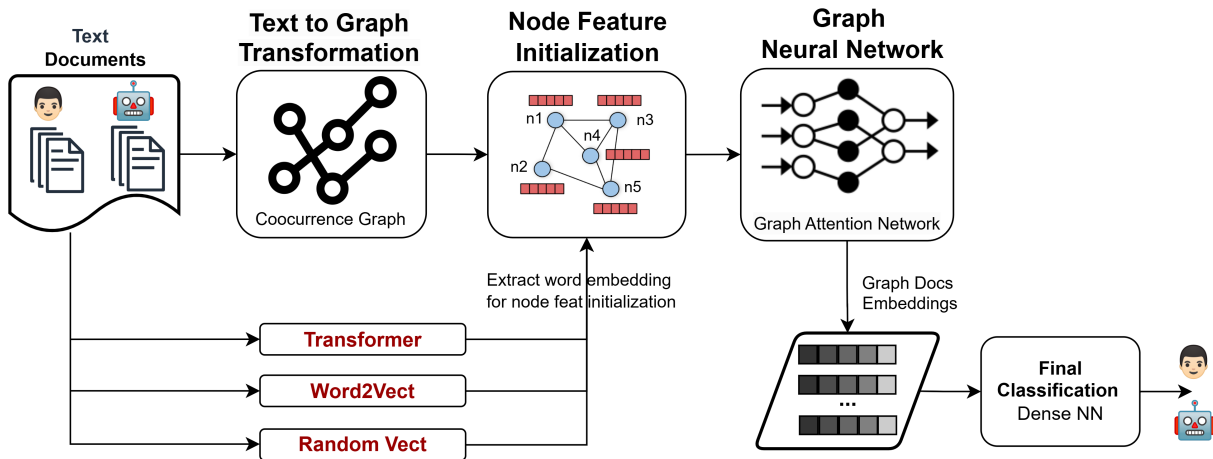


Figure 2: Pipeline Architecture for generated text identification.

Approach	Node Feat Init	Val Acc	Val F1Score	Test Acc	Test F1Score
<b>Linear-SVM</b>	-	0.7419	0.7419	0.5944	0.5624
<b>Word2Vec</b>	-	0.7325	0.7320	0.6040	0.5982
<b>FT-BERT</b>	-	0.8924	0.8916	0.6197	0.5515
<b>FT-RoBERTa</b>	-	<b>0.8974</b>	<b>0.8965</b>	0.6184	0.5481
<b>GNN Cooc-Graph</b>	Random	0.6469	0.6434	0.5738	0.5737
<b>GNN Cooc-Graph</b>	Word2Vec	0.7896	0.7896	0.6618	0.6592
<b>GNN Cooc-Graph</b>	FT-BERT	0.8889	0.8882	<b>0.7448</b>	<b>0.7441</b>
<b>GNN Cooc-Graph</b>	FT-RoBERTa	0.8812	0.8805	0.7447	0.7370

Table 2: Accuracy and Macro F1-Score on classification tasks for validation and test sets in the Autext 2023 dataset.

For the nodes feat initialization, we considered the following strategies, varying the feature vector size from 128 to 768:

- **Random.** In this approach, node features are initialized randomly using the PyTorch Embedding Layer<sup>6</sup>, taking values from -1 to 1.
- **Word2Vec.** We trained the model (on train set) using the Continuous Bag of Words method and then obtained the word embeddings for each node/word in the graph. We applied a random vector initialization for those out of the vocabulary words.
- **Transformer.** Fine-tuned the BERT and RoBERTa models using the training data and then extracted the word embeddings. The tokenizer of this model generates some fine-grained tokens for certain words. To handle and match this with the graph’s nodes, we obtain the average embedding for each token and assign the result as the node feature.

Table 2 highlights the best-performing models for text classification in the Autextification 2023 dataset, showcasing significant differences in validation and test performances across approaches. The Linear-SVM and Word2Vec baselines achieved a validation accuracy (Val F1 Macro Score) of 0.7419 and 0.7320, respectively. Among baseline models, the fine-tuned BERT and RoBERTa achieve the highest validation F1 score at 0.8965 and 0.8916, respectively, but a decline in test performance (showing high overfitting). Now, regarding the GNN Cooc-Graph, using a random feature node initialization yields moderate results, while Word2Vec improves validation and test scores (Val F1: 0.6618, Test F1: 0.6592). Using fine-tuned transformer-based features (BERT and RoBERTa models) further enhances GNN performance, achieving the best overall test results: Test F1 Score: 0.7370 for FT-RoBERTa and 0.7441 for FT-BERT. This demonstrates the effectiveness of leveraging pre-trained transformer features within a graph-based framework and how the node feature initialization significantly impacts performance.

<sup>6</sup><https://pytorch.org/docs/stable/nn.html>

Rank	Approach	Run	Macro-F1
1	TALN-UPF	HB plus	0.8091
-	GNN (our)	FT-BERT	0.7418
2	TALN-UPF	HB	0.7416
3	CIC-IPN	run2	0.7413
23	BOW+LR	baseline	0.6578
52	Transformer	baseline	0.5710
77	UAEMex	run1	0.3387

Table 3: Final Ranking on the IberAutextification 2023 shared task (English text, subtask 1).

On the other hand, Table 3 shows the final ranking on the IberAutextification 2023 shared task for subtask 1 (binary classification on English texts). The team TALN-UPF ranks 1st and 2nd with a Macro-F1 with 0.8091 and 0.7416, respectively (the system is described in section 2). As we can see, the 1st rank system outperformed our proposed approach; however, our method achieved better results than the other 76 systems, including the baselines.

## 5 Conclusions and Future work

This paper addresses the increasing prevalence of machine-generated content due to advancements in LLMs. With their rising accessibility, concerns about their misuse have grown. To tackle this, we proposed a model architecture that combines text-graph representations and GNNs to detect machine-generated text; specifically, we implemented a co-occurrence graph where each word is represented as a node, and if two words co-occur within the same text document, it is linked with an edge. Then, this graph is fed into a GNN (GAT), generating the graph document embeddings as output, which are used to train a final classification model to distinguish between human and machine text documents.

Based on the experiments in the Autextification 2023 English dataset, our approach demonstrated superior performance compared to baselines and traditional approaches, highlighting the effectiveness of integrating structural and semantic features in identifying machine-generated content.

Moreover, future studies could enhance these approaches further and investigate their applicability across various languages and datasets with different domains. Also, different text graph representations (e.g. Heterogeneous Graphs) and GNN architectures should be tried using a combination of different node feature initialization strategies.

## 6 Limitations

The results are based on a specific dataset, which may not fully represent the diversity of machine-generated content across different domains or languages. Moreover, combining GNNs and transformer-based features can be computationally expensive, making the approach less feasible for real-time and large-scale applications without further optimization.

## References

- 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*.
- Janek Bevendorff, Matti Wiegmann, Jussi Karlgren, Luise Dürlich, Evangelia Gogoulou, Arne Talman, Efsthios Stamatatos, Martin Potthast, and Benno Stein. 2024. Overview of the “Voight-Kampff” Generative AI Authorship Verification Task at PAN and ELOQUENT 2024. In *Working Notes of CLEF 2024 - Conference and Labs of the Evaluation Forum*, CEUR Workshop Proceedings. CEUR-WS.org.
- Jacob Devlin. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Yuxiao Lin, Yuxian Meng, Xiaofei Sun, Qinghong Han, Kun Kuang, Jiwei Li, and Fei Wu. 2021. Bertgcn: Transductive text classification by combining gcn and bert. *arXiv preprint arXiv:2105.05727*.
- Yinhan Liu. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 364.
- Tomas Mikolov. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 3781.
- Melania Nitu and Mihai Dascalu. 2024. [Beyond lexical boundaries: Llm-generated text detection for romanian digital libraries](#). *Future Internet*, 16(2).
- Areg Mikael Sarvazyan, José Ángel González, Marc Franco Salvador, Francisco Rangel, Berta Chulvi, and Paolo Rosso. 2023. Overview of autextification at iberlef 2023: Detection and attribution of machine-generated text in multiple domains. In *Procesamiento del Lenguaje Natural*, Jaén, Spain.
- Andric Valdez-Valenzuela and Helena Gómez-Adorno. 2024. The iimasnlp team at iberautextification 2024: Integrating graph neural networks, multilingual llms, and stylometry for automatic text identification.

- Andric Valdez-Valenzuela and Helena Gómez-Adorno. 2024. text2graphapi: A library to transform text documents into different graph representations. *SoftwareX*, 28:101888.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2017. Graph attention networks. *arXiv preprint arXiv:1710.10903*.
- Kunze Wang, Yihao Ding, and Soyeon Caren Han. 2024. Graph neural networks for text classification: A survey. *Artificial Intelligence Review*, 57(8):190.
- Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. Graph convolutional networks for text classification. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 7370–7377.