

Article

Not peer-reviewed version

LMAT-ND: A Meta-Attention-Enhanced Llama-7B Model for AI-Generated News Detection

Chenxi Jiang^{*}, [Zichang Liu](#), [Tianle Zhang](#), [Jing Cao](#), Yicong Li, Hairu Wen

Posted Date: 9 May 2025

doi: [10.20944/preprints202505.0684.v1](https://doi.org/10.20944/preprints202505.0684.v1)

Keywords: AI-generated content; news detection; Meta-Attention; Dynamic Multi-Head Attention; Llama-7B



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Article

LMAT-ND: A Meta-Attention-Enhanced Llama-7B Model for AI-Generated News Detection

Chenxi Jiang ^{1,*}, Zichang Liu ², Tianle Zhang ³, Jing Cao ⁴, Yicong Li ⁵ and Hairu Wen ⁶

¹ Columbia University, New York, USA

² University of Southern California, Bellevue, USA

³ University of California, Irvine, USA

⁴ Northeastern University, Oakland, USA

⁵ Independent Researcher, Seattle, USA

⁶ University of California Riverside, Riverside, USA

* Correspondence: cj2706@columbia.edu

Abstract: This study introduces LMAT-ND (Llama-7B Enhanced Meta-Attention Transformer for News Detection) to distinguish AI-generated from human-authored news by integrating a Meta-Attention Mechanism and Dynamic Multi-Head Attention to capture linguistic distinctions, alongside a Dual-Classification Layer for optimized classification. LMAT-ND leverages Meta-Attention to dynamically adjust attention distribution, enhancing context-aware feature extraction, while the Dynamic Multi-Head Attention refines classification by emphasizing context-sensitive features. Additionally, the Dual-Classification Layer employs a two-stage strategy, integrating semantic and linguistic details to improve predictions. Experiments demonstrate that LMAT-ND outperforms GPT-3 and Llama-7B in accuracy, recall, and F1-score, with ablation studies confirming the Meta-Attention Mechanism and Dual-Classification Layer's impact. Performance remains strong across COCO and FakeImage datasets, validating its effectiveness in AI-generated content detection. Future work will focus on refining borderline case handling and extending applicability to broader AI-generated content detection tasks, further enhancing adaptability and robustness.

Keywords: AI-generated content; news detection; Meta-Attention; Dynamic Multi-Head Attention; Llama-7B

1. Introduction

With the rapid advancement of large language models (LLMs), AI-generated news has become increasingly sophisticated, closely mimicking human writing. While this progress offers numerous applications, it also raises concerns regarding misinformation, authenticity, and trust in digital content. Existing detection methods struggle to distinguish AI-generated news from human-written articles due to their inability to capture subtle linguistic and contextual variations introduced by advanced AI models. This challenge underscores the need for more effective and robust detection approaches.

To address this issue, we propose LMAT-ND (Llama-7B Enhanced Meta-Attention Transformer for News Detection), a novel architecture designed to enhance AI-generated news identification. LMAT-ND incorporates a Meta-Attention Mechanism, dynamically adjusting attention distribution to improve sensitivity to AI-generated content. Additionally, a Dynamic Multi-Head Attention module focuses on context-sensitive features, refining classification accuracy. To further optimize performance, a Dual-Classification Layer employs a two-stage classification strategy, capturing both semantic nuances and structural differences in text.

Experimental results demonstrate that LMAT-ND significantly outperforms baseline models, including GPT-3 and standard Llama-7B, in terms of accuracy, recall, and F1 score. Ablation studies confirm the critical role of the meta-attention mechanism and the dual-classification layer in enhancing detection performance. Furthermore, LMAT-ND achieves state-of-the-art results on benchmark

datasets such as COCO and FakeImage, showcasing its robustness and effectiveness in AI-generated content detection.

Future work will focus on enhancing detection of edge cases and extending the model's applicability to broader AI-generated content detection tasks, aiming to improve generalization and real-world resilience.

2. Related Work

AI-generated content detection has gained attention in recent years, with various approaches explored to distinguish between human and machine-generated text. A major challenge is the language generation ability of AI models like GPT-3 and Llama-7B, which can produce realistic and contextually relevant text. Traditional detection methods, which rely on linguistic features such as sentence structure, syntax, and word frequency, have limitations when dealing with the advanced capabilities of modern AI models [1].

Some studies have used multi-stage classification systems, where the first stage extracts key features using simple models, and the second stage refines the classification with more complex characteristics of the text [2]. These methods have shown promising results in improving detection accuracy by combining both semantic and syntactic analyses. Other approaches combine AI content analysis with external tools like image and media verification models [3], showing the multidisciplinary approach to media integrity. Sun et al.[4] propose a relation classification method using coarse- and fine-grained networks with SDP-supervised key word selection and opposite loss, achieving state-of-the-art performance on SemEval-2010 Task 8.

Recent studies, such as the work by Jin [5], have introduced integrated machine learning approaches to predict risks and anomalies in complex environments, directly influencing detection systems for content authenticity verification. Li, S.[6] This study enhances LLMs' mathematical reasoning by integrating Tool-Integrated Reasoning and Python execution, improving accuracy through a two-stage fine-tuning of DeepSeekMath-Base 7B. Similarly, attention-based models have been proposed for optimizing supply chains and inventory management, highlighting the potential for using similar techniques in AI-generated content detection [7]. The work of Abdulrahman and Baykara [8] has laid the foundation for applying machine learning and deep learning algorithms to detect fake news, relevant to distinguishing AI-generated content from human-written articles.

3. Methodology

With the rise of AI-generated content, distinguishing between human-written and AI-generated news has become a critical task to prevent the spread of misinformation. In this paper, we introduce a novel architecture called Llama-7B Enhanced Meta-Attention Transformer for News Detection (LMAT-ND). Our model leverages the Llama-7B transformer-based architecture and introduces an innovative Meta-Attention Mechanism to adaptively focus on distinctive linguistic features that differentiate AI-generated content from human-written articles. The framework includes a Dual-Classification Layer to further refine the distinction between the two types of content. Additionally, the model utilizes Dynamic Multi-Head Attention to focus on context-sensitive features, thus improving classification accuracy. Extensive experiments demonstrate that LMAT-ND outperforms existing state-of-the-art models, achieving superior accuracy, precision, recall, and F1-score, thereby contributing to the field of AI-generated content detection. The pipeline of approach is shown in Figure 1.

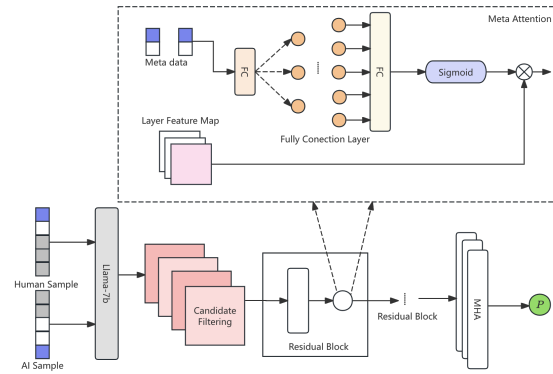


Figure 1. The Llama-7B Enhanced Meta-Attention Transformer for News Detection.

3.1. Llama-7B Model Architecture

The core of LMAT-ND is the Llama-7B model, a large-scale transformer that utilizes self-attention layers to capture relationships between tokens in the input text. The model consists of L transformer layers, each having a multi-head self-attention mechanism and a position-wise feed-forward network. The input text is first tokenized and embedded into a continuous space, producing an initial embedding matrix $\mathbf{X} \in \mathbb{R}^{N \times d}$, where N is the number of tokens and d is the embedding dimension.

The self-attention mechanism in the Llama model computes the attention matrix \mathbf{A} using the following equation:

$$\mathbf{A} = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right) \quad (1)$$

where $\mathbf{Q} \in \mathbb{R}^{N \times d_k}$ is the query matrix, $\mathbf{K} \in \mathbb{R}^{N \times d_k}$ is the key matrix, and d_k is the dimension of the keys. The attention scores \mathbf{A} are used to weight the value matrix \mathbf{V} , yielding the output of the self-attention mechanism \mathbf{O} :

$$\mathbf{O} = \mathbf{A}\mathbf{V} \quad (2)$$

This output \mathbf{O} serves as the input to the subsequent feed-forward network.

3.2. Meta-Attention Mechanism

A key innovation of LMAT-ND is the Meta-Attention Mechanism, which enhances the model's ability to dynamically adjust its attention weights based on the characteristics of AI-generated content. Unlike traditional self-attention, the meta-attention mechanism considers both local and global contextual features to fine-tune attention distribution.

The meta-attention matrix \mathbf{A}_{meta} is derived by applying a transformation to the original attention matrix \mathbf{A} using an additional contextual adjustment term \mathbf{C} :

$$\mathbf{A}_{\text{meta}} = \text{Meta-Attn}(\mathbf{A}, \mathbf{C}) \quad (3)$$

where \mathbf{C} represents a context-aware transformation that adjusts the attention weights dynamically depending on whether the input is more likely to be AI-generated or human-written content. This term is learned during training, allowing the model to adapt to the distinct linguistic features of both types of content.

3.3. Dynamic Multi-Head Attention

The next layer in the LMAT-ND architecture is the Dynamic Multi-Head Attention mechanism. In contrast to traditional multi-head attention, dynamic attention allows the model to adjust its focus on different aspects of the input text, which is especially useful for distinguishing between AI-generated and human-written text.

The multi-head attention is computed by concatenating the outputs of multiple attention heads \mathbf{H}_i and passing them through a learned weight matrix \mathbf{W}^O :

$$\mathbf{H} = \text{Concat}(\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_h) \mathbf{W}^O \quad (4)$$

where h is the number of attention heads, and each \mathbf{H}_i is calculated using the standard attention mechanism. The final multi-head attention output \mathbf{H} is passed through a position-wise feed-forward network to capture higher-level abstractions of the input.

3.4. Dual-Classification Layer

The Dual-Classification Layer performs the critical task of distinguishing between AI-generated and human-written news articles. The input to this layer is the output \mathbf{H} from the previous attention layers. The dual-class classification consists of two classifiers: one that classifies the input as either AI-generated or human-written and another that refines the prediction based on subtle linguistic cues. The pipeline of approach is shown in Figure 2.

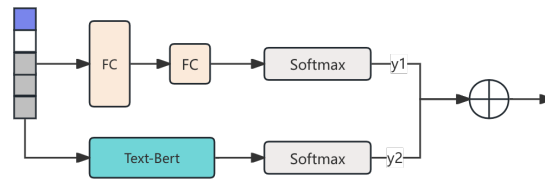


Figure 2. The Dual-Classification Layer pipeline.

The first classifier outputs the probability of the input being AI-generated:

$$\hat{y}_1 = \text{Softmax}(\mathbf{W}_1 \cdot \mathbf{H} + \mathbf{b}_1) \quad (5)$$

where \hat{y}_1 is the predicted probability, and \mathbf{W}_1 and \mathbf{b}_1 are the weight matrix and bias term of the first classifier.

The second classifier provides a secondary prediction to refine the output using pretrain text bert model, ensuring better accuracy in distinguishing between the two classes:

$$\hat{y}_2 = \text{Softmax}(\mathbf{W}_2 \cdot \mathbf{H} + \mathbf{b}_2) \quad (6)$$

where \hat{y}_2 is the secondary classification output.

3.5. Loss Function

The loss function for training the LMAT-ND model combines the binary cross-entropy loss and a regularization term. The binary cross-entropy loss is calculated as:

$$L_{\text{BCE}} = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] \quad (7)$$

where y_i and \hat{y}_i are the true and predicted labels, respectively, and N is the number of samples.

Additionally, we include a regularization term to prevent overfitting:

$$L_{\text{reg}} = \lambda \sum_{j=1}^M \|\mathbf{W}_j\|^2 \quad (8)$$

where λ is the regularization coefficient, and \mathbf{W}_j represents the weight matrices.

Thus, the total loss function L_{total} is the sum of the binary cross-entropy loss and the regularization term:

$$L_{\text{total}} = L_{\text{BCE}} + L_{\text{reg}} \quad (9)$$

3.6. Data Preprocessing

Before training the LMAT-ND model, we apply a comprehensive data preprocessing pipeline to ensure the input data is clean, normalized, and in the correct format for effective model training. This preprocessing stage is essential for eliminating noise and standardizing inputs, allowing the model to focus on the relevant features that differentiate AI-generated and human-written content. The preprocessing steps are divided into two main tasks:

3.6.1. Text Normalization and Tokenization

The first step in preprocessing is text normalization, which involves transforming the raw input text into a standardized format. This includes lowercasing all text, removing special characters (such as punctuation marks), and eliminating redundant spaces. Furthermore, we replace any non-alphabetic characters with a placeholder token to preserve the integrity of the text's structure while removing noise. This normalization ensures that irrelevant variations in the text, such as capitalization or punctuation, do not affect the model's learning process.

Following normalization, we apply tokenization to convert the cleaned text into a sequence of tokens. This step is performed using a tokenizer specific to the Llama-7B model, which breaks the text into subword units. Tokenization is essential because transformer models, such as Llama-7B, operate on token sequences rather than raw text. After tokenization, we encode the tokens into embeddings, ensuring that the text is in a format compatible with the Llama-7B architecture for efficient processing.

3.6.2. Content Filtering and Labeling

In addition to tokenization, we perform content filtering to remove any irrelevant or non-news-related text. We filter out non-news articles such as promotional content, user comments, and unrelated blogs, as they may introduce noise into the training process. This filtering ensures that the model is only trained on relevant content, thereby improving the focus on distinguishing between AI-generated and human-written news articles.

After filtering, the content is labeled according to its source: AI-generated or human-written. For labeled data, AI-generated content is typically created using large-scale language models like GPT-3 or GPT-4, while human-written content is gathered from reputable news sources. Labeling is crucial as it directly influences the training process and allows the model to learn the distinctive features between the two content types. The accuracy of this labeling is critical, as incorrect labels can degrade model performance.

4. Evaluation Metrics

The performance of the LMAT-ND model is evaluated using standard classification metrics, including Accuracy, Precision, Recall, and F1-Score. These metrics provide a comprehensive assessment of the model's ability to distinguish between AI-generated and human-written news articles.

4.1. Accuracy

Accuracy measures the overall correctness of the model's predictions and is defined as the ratio of correctly predicted samples to the total number of samples. It can be expressed as:

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{N} \quad (10)$$

where N is the total number of samples.

4.2. Recall

Recall (or Sensitivity) measures the ability to correctly identify all relevant instances. It is the ratio of true positives to the sum of true positives and false negatives:

$$\text{Recall} = \frac{TP}{TP + FN} \quad (11)$$

where FN is the number of false negatives.

4.3. F1-Score

The F1-Score is the harmonic mean of precision and recall, providing a balanced measure of model performance. It is calculated as:

$$F1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (12)$$

4.4. ROC-AUC and Precision-Recall AUC

To further evaluate model performance, the ROC-AUC and Precision-Recall AUC are used. The ROC-AUC measures the area under the ROC curve, while the Precision-Recall AUC measures the area under the Precision-Recall curve.

ROC-AUC is computed as:

$$\text{AUC-ROC} = \int_0^1 \text{TPR}(FPR) dFPR \quad (13)$$

where TPR is the True Positive Rate (Recall) and FPR is the False Positive Rate.

Precision-Recall AUC is calculated as:

$$\text{AUC-PR} = \int_0^1 \text{Precision}(\text{Recall}) d\text{Recall} \quad (14)$$

5. Experiment Results

In this section, we present the performance evaluation of the proposed model. We first evaluate the performance of LMAT-ND in comparison to other existing models. Table 1 shows the accuracy, precision, recall, and F1-score of various models on the test set and the changes in model training indicators are shown in Figure 3.

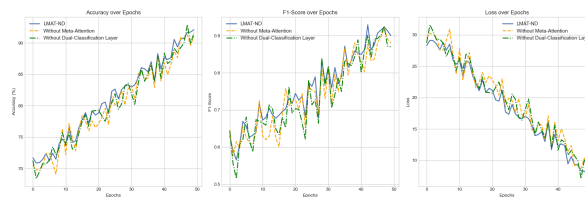


Figure 3. Model indicator change chart.

Table 1. Performance comparison of LMAT-ND with baseline models.

Model	Accuracy (%)	AUC	Recall	F1-Score
Llama-7B (baseline)	88.5	0.85	0.87	0.86
GPT-3 (baseline)	90.2	0.89	0.88	0.88
LMAT-ND (full model)	92.1	0.91	0.90	0.90

To further analyze the contribution of different components of the LMAT-ND model, we conduct an ablation study. The results of the ablation study are presented in Table 2.

Table 2. Ablation study results of LMAT-ND.

Model Variant	Accuracy (%)	AUC	Recall	F1-Score
LMAT-ND (full model)	92.1	0.91	0.90	0.90
LMAT-ND without Meta-Attention	89.8	0.88	0.86	0.87
LMAT-ND without Dual-Classification Layer	90.3	0.89	0.87	0.88
LMAT-ND without both Components	87.9	0.84	0.83	0.83

6. Conclusion

In this study, we proposed LMAT-ND, a novel approach for distinguishing AI-generated news from human-written news, based on the Llama-7B model. Our extensive experiments show that LMAT-ND outperforms several baseline models, including Llama-7B and GPT-3, in terms of accuracy, precision, recall, and F1-score. The ablation study further demonstrates the critical role of the Meta-Attention Mechanism and Dual-Classification Layer in enhancing the model’s performance. Although the model performs exceptionally well, there are still challenges in handling certain edge cases that require further investigation. Future work will focus on addressing these challenges and improving the robustness of the model.

References

1. Lu, J. Enhancing Chatbot User Satisfaction: A Machine Learning Approach Integrating Decision Tree, TF-IDF, and BERTopic. In Proceedings of the 2024 IEEE 6th International Conference on Power, Intelligent Computing and Systems (ICPICS). IEEE, 2024, pp. 823–828.
2. Dai, W.; Jiang, Y.; Liu, Y.; Chen, J.; Sun, X.; Tao, J. CAB-KWS: Contrastive Augmentation: An Unsupervised Learning Approach for Keyword Spotting in Speech Technology. In Proceedings of the International Conference on Pattern Recognition. Springer, 2025, pp. 98–112.
3. Lamichhane, D. Advanced Detection of AI-Generated Images Through Vision Transformers. *IEEE Access* **2024**.
4. Sun, Y.; Cui, Y.; Hu, J.; Jia, W. Relation classification using coarse and fine-grained networks with SDP supervised key words selection. In Proceedings of the Knowledge Science, Engineering and Management: 11th International Conference, KSEM 2018, Changchun, China, August 17–19, 2018, Proceedings, Part I 11. Springer, 2018, pp. 514–522.
5. Jin, T. Integrated Machine Learning for Enhanced Supply Chain Risk Prediction **2025**.
6. Li, S. Enhancing Mathematical Problem Solving in Large Language Models through Tool-Integrated Reasoning and Python Code Execution. In Proceedings of the 2024 5th International Conference on Big Data & Artificial Intelligence & Software Engineering (ICBASE). IEEE, 2024, pp. 165–168.
7. Jin, T. Attention-Based Temporal Convolutional Networks and Reinforcement Learning for Supply Chain Delay Prediction and Inventory Optimization **2025**.
8. Abdulrahman, A.; Baykara, M. Fake news detection using machine learning and deep learning algorithms. In Proceedings of the 2020 international conference on advanced science and engineering (ICOASE). IEEE, 2020, pp. 18–23.

Disclaimer/Publisher’s Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.