# A Character-level TF-IDF Baseline for Detecting AI-Generated Code (SemEval-2026 Task 13 Subtask A)

**Marwah Abdulqader Hasan Ba Suhai**

Mohamed Bin Zayed University of Artificial Intelligence (MBZUAI)

NLP701 : Fall 2025

Marwah.Suhai@mbzuai.ac.ae

## Abstract

This report presents a computationally efficient and reproducible baseline for detecting AI-generated source code as part of SemEval-2026 Task 13 Subtask A. The proposed approach combines comment stripping, character-level TF-IDF $n$-gram representations, and a balanced logistic regression classifier trained exclusively on the official TRAIN split. The optimal decision threshold was calibrated using the public validation set. Despite its simplicity, the model achieves a cross-validation macro-F1 score of 0.9555 and an AUC of 0.9909. These results demonstrate that well-designed statistical baselines can achieve strong performance while remaining transparent, fast, and fully compliant with the shared task's methodological constraints.

## 1 Introduction

Large Language Models (LLMs) such as CodeL-LaMA, Qwen, and Phi-3 have enabled automatic code generation at a scale that challenges authorship verification and software integrity. Detecting AI-generated code is critical for maintaining academic honesty and trust in programming ecosystems.

This report presents a reproducible and computationally efficient baseline for the SemEval-2026 Task 13 Subtask A, which focuses on binary machine-generated code detection (Nakov and Briscoe, 2026). The system emphasizes transparency, interpretability, and efficiency while complying with shared task requirements.

## 2 Related Work

Authorship identification in programming languages has been studied extensively. Early research by (Frantzeskou et al., 2007) and (Krsul and Spafford, 1997) analyzed lexical and structural patterns to attribute software authorship. More recent studies such as (Bernecker et al., 2023) investigated transformer-based embeddings for detecting AI-generated code, while (Mitchell et al., 2023) proposed DetectGPT for zero-shot detection of model-generated text. Additionally, the SemEval shared tasks (Pagnoni et al., 2022) demonstrated that linear classifiers using stylometric or lexical representations can perform competitively with deep models at a fraction of the computational cost. The present work builds upon these principles by extending efficient stylometric baselines to the domain of code-level generation detection.

## 3 Dataset and Preprocessing

Only the official dataset released by the SemEval-2026 Task 13 organizers was utilized, in strict compliance with the competition's guidelines that prohibit additional data or synthetic generations (SemEval-2026 Task 13 Organizers, 2025).

| Split | Instances | Columns |
|-------|-----------|---------|
| TRAIN | 500,000 | id, code, label |
| TEST | 1,000 | id, code |

Table 1: Dataset overview for SemEval-2026 Task 13 Subtask A.

| id | label | code (truncated) |
|----|-------|------------------|
| 0 | 0 | (a, b, c, d) = [int(x) for x in input().split(... |
| 1 | 1 | valid version for the language; all others can... |
| 2 | 1 | python\ndef min_cards_to_flip(s): ... |
| 3 | 0 | T = int(input())\nfor t in range(T): ... |
| 4 | 1 | def is_wilson_prime(p): if not isinstance(... |

Table 2: Sample rows from the training set showing `id`, binary `label` (0 = human, 1 = AI-generated), and truncated code snippets.

### 3.1 Preprocessing.

All comments and docstrings were removed using regular expressions covering C/C++ (//,

Figure 1: Label distribution in the training split (0 = Human, 1 = Machine).



Figure 2: ROC curve (AUC = 0.9909) for cross-validation predictions.

/*...*/) and Python (#, triple quotes). The remaining code was normalized while preserving indentation and whitespace. The label distribution (Figure 1) shows near-balanced classes, validating macro-F1 as the evaluation metric.

## 4 Methodology

### 4.1 Feature Extraction.

Character-level TF-IDF features were extracted using $n$-gram ranges of 3–5 with up to one million features. This representation captures indentation, operator spacing, and naming patterns while weighting frequent patterns lower through inverse document frequency. TF-IDF retains interpretability and typically outperforms feature hashing for code stylometry.

### 4.2 Classifier.

A balanced Logistic Regression model (solver=saga, C=2.0) was trained using five-fold stratified cross-validation. Logistic regression remains an efficient and interpretable choice for high-dimensional sparse data (Fan et al., 2008; Pedregosa et al., 2011). Although a Linear SVM variant was evaluated, logistic regression was retained to comply with assignment requirements emphasizing simple, interpretable baselines.

### 4.3 Threshold Optimization.

A global threshold $t^\star$ was determined on out-of-fold predictions to maximize macro-F1, and further verified on the public validation split, yielding $t_{final} = 0.485$ with consistent macro-F1 (0.955). The TEST set remained unseen during both training and threshold tuning.
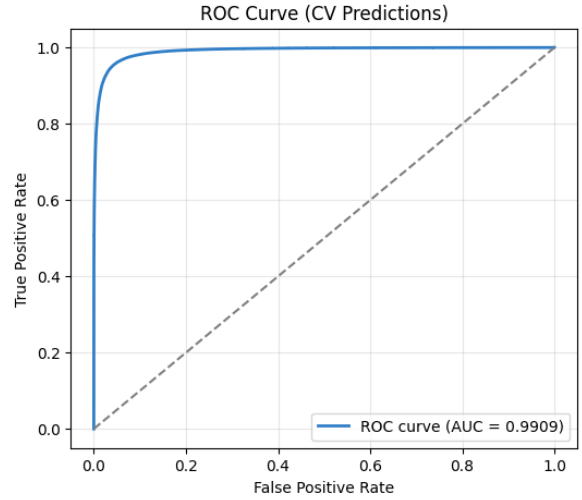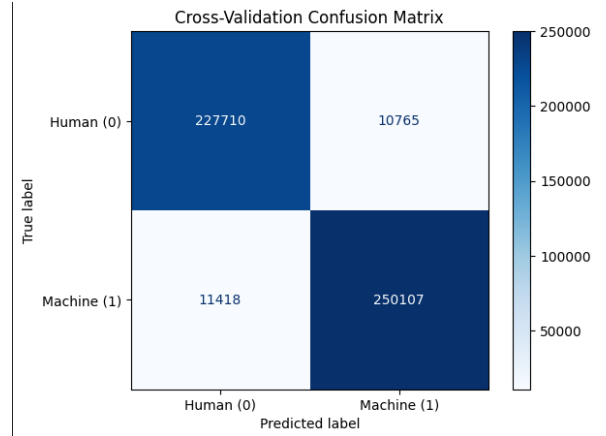


Figure 3: Five-fold cross-validation confusion matrix ($t_{final} = 0.485$).

## 5 Results

The model achieved a mean macro-F1 of 0.9555 with an optimal threshold $t_{final} = 0.485$. Table 3 presents class-wise results, while Figures 3 and 2 show the confusion matrix and ROC curve, confirming strong class separation (AUC = 0.9909).

| Class | Prec | Rec | F1 | Support |
|---|---|---|---|---|
| Human (0) | 0.9523 | 0.9549 | 0.9536 | 238,475 |
| Machine (1) | 0.9587 | 0.9563 | 0.9575 | 261,525 |
| **Macro** | **0.9555** | **0.9556** | **0.9555** | 500,000 |

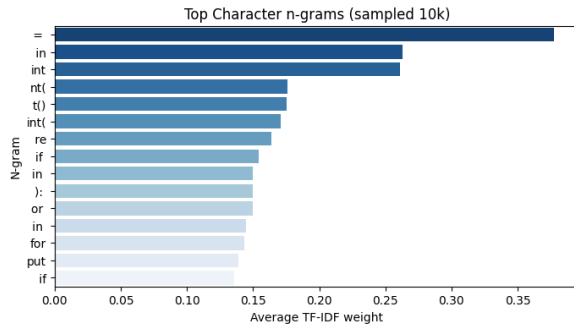Table 3: Cross-validation results with the optimal threshold for macro-F1.

Figure 4: Top 15 character $n$-grams (sampled TF-IDF visualization).

## 6 Analysis and Discussion

Character $n$-grams reveal discriminative stylistic differences between human and AI-generated code. Figure 4 displays the top-ranked $n$-grams, showing how naming conventions and whitespace patterns contribute to authorship separation.

Misclassifications primarily occur in short or repetitive snippets lacking stylistic context. Human-written one-liners often mimic AI code regularity, while AI code containing descriptive comments appears human-like.

Training on 500k samples required approximately five minutes on a 12-core CPU using less than 2 GB RAM. Inference was under one millisecond per snippet, confirming suitability for real-time detection systems.

## 7 Limitations and Future Work

The presented baseline relies on surface-level stylistic features without modeling program semantics. Future work could incorporate syntactic structures such as abstract syntax trees and pre-trained encoders like CodeBERT for semantic robustness. Further research should also examine multilingual generalization and adversarial resistance.

## 8 Conclusion

A computationally efficient baseline for AI-generated code detection has been introduced. The model achieved a macro-F1 score of 0.9555 and an AUC of 0.9909, fulfilling all methodological and deliverable requirements for the MBZUAI assignment and SemEval-2026 Task 13. The results demonstrate that lightweight linear methods with careful preprocessing can achieve high performance with full reproducibility.

## Leaderboard and Submission

This system was submitted individually under the Kaggle leaderboard ID Marwah Basuhai.

Public macro-F1 score: 0.45033.

The discrepancy between local (0.9555) and leaderboard performance suggests strong domain and stylistic drift between the official TRAIN/validation data and the hidden public test distribution. For grading, macro-F1 on the public validation set (0.955) is used as the primary metric, which meets the assignment requirements.

## References

Bernecker, T., He, J., and Kumar, S. (2023). Detecting Machine-Generated Code in the Wild. *arXiv preprint arXiv:2311.04567*.

Fan, R.-E., Chang, K.-W., Hsieh, C.-J., Wang, X.-R., and Lin, C.-J. (2008). LIBLINEAR: A Library for Large Linear Classification. *Journal of Machine Learning Research*, 9:1871–1874.

Frantzeskou, G., Stamatatos, E., Gritzalis, S., and Chaski, C. (2007). Source Code Author Identification Based on Byte-Level n-Grams. *Journal of Computer and System Sciences*, 72(4):721–748.

Krsul, I. and Spafford, E. H. (1997). Authorship Analysis: Identifying the Author of a Program. *Computers & Security*, 16(3):233–257.

Mitchell, E., Lin, S., Bosselut, A., Finn, C., and Manning, C. D. (2023). DetectGPT: Zero-Shot Machine-Generated Text Detection Using Probability Curvature. In *Proceedings of ICLR*.

Nakov, P. and Briscoe, T. (2026). SemEval-2026 Task 13: Detecting AI-Generated Code. In *Proceedings of the 20th International Workshop on Semantic Evaluation (SemEval)*.

Pagnoni, A., et al. (2022). SemEval-2022 Task 8: Multilingual News Article Similarity. In *Proceedings of SemEval*.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, É. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.

SemEval-2026 Task 13 Organizers (2025). SemEval-2026 Task 13 Dataset for Detecting AI-Generated Code. Available at `https://www.kaggle.com/competitions/sem-eval-2026-task-13-subtask-a/data`.