

CLEF 2025 PAN: Multi-Author Writing Style Analysis

Ali Muhammad Asad¹, Musaib², Syed Muhammad Areeb Kazmi³ and Sarim Tahir⁴

¹Department of Computer Science, Habib University

²Department of Computer Science, Habib University

³Department of Computer Science, Habib University

⁴Department of Computer Science, Habib University

Abstract

Multi-Author Writing Style Analysis is a crucial task in computational linguistics and authorship attribution, aimed at identifying points of transition within a document. This task has significant applications in forensic linguistics, plagiarism detection, and content verification. Over the years, various approaches have been explored, ranging from traditional stylometric techniques to advanced deep learning models. Early research relied on lexical and syntactic features to detect style changes, but these methods faced limitations in handling nuanced transitions. Recent advancements, particularly the integration of transformer-based models, have significantly improved the accuracy of style change detection. This paper aims to provide an overview of the tasks, the datasets, the evaluation metrics, and the baseline models for the Multi-Author Writing Style Analysis Lab of the PAN track at CLEF 2025.

1. Introduction

Writing style analysis is a fundamental problem in computational linguistics, with applications in authorship attribution, forensic text analysis, and collaborative content verification. A challenge in this domain is detecting when authorship changes within a multi-author document, a task known as style change detection. This problem has been the focus of multiple iterations of the PAN shared tasks, which have provided benchmark datasets and evaluation frameworks to advance the field.

Early approaches to style change detection leveraged handcrafted features such as n-grams, part-of-speech tag frequencies, and sentence structure analysis. These methods, while effective for basic segmentation tasks, struggled with complex, subtle changes in authorial style. The introduction of deep learning and pre-trained transformer models revolutionized the field, allowing for more nuanced analysis of writing styles. Models like BERT and DeBERTa have demonstrated superior performance in detecting stylistic shifts at both the paragraph and sentence levels.

This paper aims to explore existing methods and challenges in multi-author writing style analysis, reviewing state-of-the-art approaches such as deep learning models, transfer learning, and ensemble techniques. We will discuss the datasets, evaluation metrics, and baseline models used in the previous iterations of Multi-Author Writing Style Analysis Lab of the PAN track, ultimately attempting to improve the performance of style change detection.

2. Literature Review

Traditional approaches to multi-author writing style analysis have relied on lexical and syntactic feature-based methods, including TF-IDF vectorization, character n-grams, and POS tag frequencies. These techniques have demonstrated moderate performance in detecting style changes at the paragraph level but struggle with sentence-level changes due to limited contextual understanding. One study leveraged a combination of lexical and structural features, such as indentation and sentence length, to refine style change detection by incorporating discourse-level analysis. However, handcrafted features

were constrained by their susceptibility to topic influence, leading to inconsistent performance across datasets [1] [2].

More recent approaches of 2024 and 2023 have shifted towards deep learning and transformer-based models to address these limitations. A paper in PAN 2024 explored bagging techniques, feature engineering, BERT-based classifiers and ensemble techniques combining different architectures such as BERT, RoBERTa, Electra and Llama2, using DetectGPT (Mistral-7B Falcon-7B) as their baseline models. They were able to achieve an F1-score of 0.924 as their best results on fine-tuned Mistral and Llama2 models [3]. Another noteworthy paper of 2023 leveraged pre-trained transformers like DeBERTaV3 and BERT for paragraph-level classification, achieving state-of-the-art F1 scores of 0.83 for the hardest dataset variant [4]. These models incorporated contrastive learning and fine-tuning on task-specific datasets, allowing them to distinguish between subtle stylistic variations more effectively. However, their performance degraded when topic variations were minimized, highlighting the challenge of isolating style from content.

Another approach was used in 2023 where the authors performed data augmentation before training several models including RoBERTa, ELECTRA, and BERT. They achieved the best scores on RoBERTa with F1 scores of 0.996, 0.811, 0.814 for easy, medium, and hard datasets respectively [5]. Another paper decided to use a different approach where they used a contrastive learning method to optimize the segment embedding output by the encoder of the pre-training model to obtain more similar vector spaces when processing sentences with similar styles. They called this CoSENT, and were able to get F1 scores of 0.915, 0.820, and 0.705 on easy, medium, and hard datasets respectively [6].

Tzu-Mi Lin et.al experimented with ensemble pre-trained transformer models including BERT, RoBERTa, and ALBERT. The best scores they got were for RoBERTa overall with F1 scores of 0.766, 0.503, and 0.705 for easy, medium, and hard datasets respectively [7]. Hybrid Deep learning models have also been explored such as Bi-LSTM and BERT where BERT was used to create word embeddings which were then passed onto a BI-LSTM layer, followed by a CNN layer, then by a pooling layer, and finally through a fully connected layer [8]. This model achieved F1 scores of 0.67, 0.40, and 0.65 for easy, medium, and hard datasets respectively.

3. Methodology

3.1. Dataset

The dataset was provided by PAN, based on user posts from various subreddits of the Reddit platform. It is divided into three levels: easy, medium and hard, where each level is split into three parts:

- *training set*: Contains 70% of the whole dataset and includes ground truth data. This data would be used to develop and train the models.
- *validation set*: Contains 15% of the whole dataset and includes ground truth data. This data would be used to evaluate and optimize the models.
- *test set*: Contains 15% of the whole dataset and does not include ground truth data. This data would be used to evaluate the models.

Input Format

For each problem instance X (i.e., each input document), two files are provided:

1. *problem-X.txt* which contains the actual text in the form of sentences of varying lengths.
2. *truth-problem-X.json* which contains the ground truth, i.e., the correct solution in JSON format.

A sample json file looks as so:

```
{"authors": 2,  
"changes": [0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0]}
```

where the key 'changes' is an array of consecutive sentences within each document (0 when there is no change, and 1 when there is a change).

3.2. Initial Approach

We used the baseline provided in PAN 2024 Overview [3] for comparison for our models and approaches. Our initial approach was to use the Bag of Words (BoW) method for feature extraction from the sentences, coupled with a simple logistic regression model for classification. The Bag of Words method is a simple and effective way to represent text data. It involves creating a vocabulary of unique words and then representing each sentence as a vector of word counts. This became our baseline approach, as it is simple and easy to implement, and was backed by previous implementations [3].

Then we modified the bag of words approach to combine n-grams and syntactic/lexical features (sentence length, POS tag frequencies, etc), along with calss weighting to handle the class imbalance between instances of 0s and 1s, and finally used a vectorized tuning with ngram_range of (1, 2) to include both unigrams and bigrams in the feature set, in order to improve the model's performance by capturing more context and semantics from the text data. This showed an improvement in the model's performance over the baseline.

4. Results

This section presents the results of our experiments. The table below shows the performance of the models on the validation set. The models were evaluated based on the F1 score, which is the harmonic mean of precision and recall. The F1 score is a good metric for imbalanced datasets, as it considers both false positives and false negatives. The F1 score is presented for each level, easy, medium and hard.

Level	F1 (Easy)	F1 (Medium)	F1 (Hard)
Baseline (2024 Results)	0.414	0.506	0.495
Simple BoW	0.653	0.650	0.573
Improved BoW	0.791	0.666	0.613

Table 1

F1 scores of the models on the validation set

The results show an improvement in the F1 score for the improved BoW model over the baseline and BoW model on all levels. We will next further improve upon our results by exploring ML models along with TF-IDF vectorization as they can better capture the semantics of the text, we will also train Deep learning models such as Bi-LSTM, and transformer based models such as DeBertaV3, RoBERTa, BERT, Mistral, and LLama models as they have been proved by literature to perform well in the previous iterations of this task.

References

- [1] E. Zangerle, M. Mayerl, M. Potthast, B. Stein, Overview of the style change detection task at pan 2021, CLEF 2021 - Conference and Labs of the Evaluation Forum, September 21-24, 2021, Bucharest, Romania (2021).
- [2] E. Zangerle, M. Mayerl, M. Potthast, B. Stein, Overview of the style change detection task at pan 2022, CLEF 2022 - Conference and Labs of the Evaluation Forum, September 5-8, 2022, Bologna, Italy (2022).
- [3] A. A. Ayele, N. Babakov, J. Bevendorff, X. B. Casals, B. Chulvi, D. Dementieva, A. Elnagar, D. Freitag, M. Fröbe, D. Korencic, M. Mayerl, D. Moskovskiy, A. Mukherjee, A. Panchenko, M. Potthast, F. Rangel, N. Rizwan, P. Rosso, F. Schneider, A. Smirnova, E. Stamatatos, E. Stakovskii, B. Stein, M. Taulé, D. Ustalov, X. Wang, M. Wiegmann, S. M. Yimam, E. Zangerle, Overview of pan 2024: Multi-author writing style analysis, multilingual text detoxification, oppositional thinking analysis, and generative ai authorship verification condensed lab overview, CLEF 2024: Conference and Labs of the Evaluation Forum, September 09-12, 2024, Grenoble, France (2024).

- [4] E. Zangerle, M. Mayerl, M. Potthast, B. Stein, Overview of the multi-author writing style analysis task at pan 2023, CLEF 2023 - Conference and Labs of the Evaluation Forum, September 18-21, 2023, Thessaloniki, Greece (2023).
- [5] A. Hashemi, W. Shi, Enhancing writing style change detection using transformer-based models and data augmentation, CLEF 2023 - Conference and Labs of the Evaluation Forum, September 18-21, 2023 (2023).
- [6] H. Chen, Z. Han, Z. Li, Y. Han, A writing style embedding based on contrastive learning for multi-author writing style analysis, CLEF 2023 - Conference and Labs of the Evaluation Forum, September 18-21, 2023 (2023).
- [7] T.-M. Lin, C.-Y. Chen, Y.-W. Tzeng, L.-H. Lee, Ensemble pre-trained transformer models for writing style change detection, Conference and Labs of the Evaluation Forum, CLEF 2022 (2022).
- [8] J. Zi, L. Zhou, Z. Liu, Style change detection based on bi-lstm and bert, Conference and Labs of the Evaluation Forum, CLEF 2022 (2022).