

# MBTI Personality Classification through Analyzing CVs/ Personal Statements for e-Recruitment

Himaya Perera\* and Lakshan Costa

Himaya Perera and Lakshan Costa are with the Informatics Institute of Technology, Sri Lanka.

Email: himaya.2019379@iit.ac.lk.org (H.P.); lakshan.c@iit.ac.lk(L.C.)

\*Corresponding author

Manuscript received May 10, 2023; revised June 14, 2023; accepted July 2, 2023; published March 15, 2024

**Abstract**—With the advent of the internet and globalization, most services have become automated, including the recruitment process. E-recruitment systems have gained popularity due to their ability to automate various steps of recruitment. However, while some systems filter and screen CVs for skills, they do not assess a candidate's personality, which is crucial in determining their compatibility with a company's culture and practices. This compatibility leads to effective recruitment and happier, more productive employees who stay at the company longer. Due to the lack of a personality detection mechanism, recruiters spend a significant amount of time conducting multiple interviews to assess candidate fit. To address this issue, the authors propose a system that uses fine-tuned BERT to detect a candidate's MBTI personality from their CV or personal statement. The proposed system achieved an average accuracy of 72.14% per class and an AUC range of 0.85-0.90, following pre-processing of the Kaggle MBTI personality types dataset to address class imbalance.

**Keywords**—MBTI, personality, multi-class classification

## I. INTRODUCTION

Over time, e-recruitment systems have become the primary means of recruitment for many large companies and businesses worldwide. This shift has been fueled by the rise of social media platforms like LinkedIn, as well as globalization and widespread internet access. Unlike the traditional process of mailing CVs to companies in hopes of securing an interview, candidates are now required to upload their CV or provide their details to e-recruitment sites where matching candidates to jobs occurs more seamlessly [1]. With the increasing prevalence of machine learning (ML) and artificial intelligence (AI) technologies, automation has become a vital component of many industries, making processes faster and more efficient while reducing manual labour. In light of the demanding nature of recruitment, the author argues that automating the hiring process through the matching of candidates to jobs based not only on skill but also on personality will greatly enhance its value and effectiveness.

The majority of e-recruitment systems have made significant strides towards proficiency in scouring candidates' CVs for requisite skills relevant to a particular job. These searches employ varying degrees of technological sophistication, ranging from rudimentary word-matching algorithms to more comprehensive assessments of job compatibility, which take into account an individual's search history and past employment experiences, reminiscent of popular social networking platforms such as LinkedIn. However, despite the evolution of these search capabilities, there remains a critical limitation in terms of

mapping candidates to roles that fit the company's core values and assessing their adaptability to the given work environment based on their unique personality traits.

Personality, as defined by scholars, represents the core characteristics of an individual that determine their perceptions, emotions, and actions [2]. Each person possesses a distinctive configuration of traits that shapes their behavior as an individual. A candidate's personality is a critical determinant of their job performance when employed by a specific organization. The degree of compatibility between a candidate and a company is a reliable indicator of their level of satisfaction and motivation while working, and ultimately, an accurate predictor of their output [3]. Research suggests that relying solely on past work experience to predict job performance only yields an accuracy rate of 16%. In contrast, the assessment of a candidate's personality traits can predict job performance with a much higher accuracy rate of approximately 78% [4].

The problem statement is defined as follows: The current state of e-recruitment systems is such that they do not take into account a candidate's personality in addition to their required skills when matching them to jobs. This presents a problem, and in this project, the author aims to address it by exploring the following research questions: What is the importance of personality when evaluating a candidate for a job? What technology can be utilized to predict personality based on collected data? How can an individual's MBTI type be calculated from textual inputs? And should personality type be classified using the four MBTI dichotomies? The desired outcome of this project is a program that accurately identifies an individual's MBTI personality type within the four dichotomies of MBTI.

This paper provides the following contributions. (1) comparison of various approaches for this multiclass classification. (2) The Dataset preprocessing. (3) training, finetuning and evaluation of BERT for MBTI classification.

The structure of this paper is outlined as follows: Section II breaks down the concept of MBTI personality indicator, specifically compared to the Big Five characteristics. Section III presents similar and related work. Section IV presents the methods of data preprocessing and training and finetuning the model. Section V presents the test results and evaluates the usefulness of the proposed method in the selected domain of e-recruitment.

## II. PERSONALITY CLASSIFICATIONS AND THE IMPORTANCE OF PERSONALITY FOR RECRUITMENT

### A. The Big Five

Psychologists have made various attempts to determine

the precise number of traits that make up an individual's personality. Historically, estimates of the number of traits have varied widely, ranging from 40,000 as suggested by Gordon Allport, to as few as 3, as proposed by Hans Eysenck, with Raymond Cattell suggesting 16 traits [5]. However, many psychologists found these estimates to be either too high or too low, which led to the emergence of "The Big Five" model. This model, which includes openness, conscientiousness, extraversion, agreeableness, and neuroticism, was later found to account for the data exceptionally well, as demonstrated by researchers such as Tupes and Christal who reanalyzed the same data used by other psychologists [6].

The Big Five are namely,

- 1) Extraversion: the degree to which someone is sociable or outgoing. Those on the high end are extroverts, while those on the lower end are introverts who may need more alone time.
- 2) Agreeableness: the level of kindness, trust, and cooperativeness. High scorers tend to be cooperative, while those on the lower end may be more competitive.
- 3) Openness: the level of creativity and interest in new experiences. Those on the extreme end tend to be adventurous and creative, while those on the lower end are more traditional and struggle with creativity.
- 4) Conscientiousness: the level of thoughtfulness, impulse control, and goal orientation. High scorers are typically organized and detail-oriented, while those on the lower end may be more scattered and less organized.
- 5) Neuroticism: the level of emotional instability and moodiness. High scorers may experience anxiety and sadness, while those on the lower end tend to be more emotionally stable and resilient [5].

Most researchers generally agree with the concept of the Big Five, but there are some minor disagreements about how to categorize certain traits. The primary issue is whether to divide the dimension of "extraversion" into the dimensions "ambition" and "sociability".

### B. The Myer Briggs Type Indicator

Isabel Briggs Myers and her mother Katharine Briggs developed the Myers-Briggs Type Indicator (MBTI) personality, which is based on C.G. Jung's theory of psychological types [7]. The MBTI personality consists of 16 personality types, which are distinguished by the following four dichotomies:

- 1) Extraversion or Introversion (E or I): Do individuals prefer the external or internal world?
- 2) Sensing or Intuition (S or N): Do individuals focus on basic information or interpret and add meaning to it?
- 3) Thinking or Feeling (T or F): Do individuals prioritize logic or consider the people and circumstances involved in decision-making?
- 4) Judging or Perceiving (J or P): Do individuals prefer to have things decided or remain open to new information and perspectives?

Each dichotomy is represented by one letter when describing an individual's MBTI personality type [7]. For instance, an ESTJ personality type is an extraverted, sensitive individual who thinks rather than feels and is judgmental. Considering all of the above dichotomies results

in 16 MBTI personality types. This project utilizes MBTI personality types to describe a user's personality as they provide a more detailed and nuanced description compared to simply identifying an individual's Big Five types.

### C. The Importance of Personality for Recruitment

Assessing a candidate's personality and background is vital for an organization as it is a contributing factor to an employee's performance. Personality often plays a part in how well an individual fits in at a company and can adapt and perform at their best within the workplace environment [6]. Hypotheses that conscientiousness and low neuroticism are valid predictors of job performance across all job types since individuals who are more thoughtful and emotionally stable can be more effective at their jobs. Similarly, individuals with high neuroticism will be more emotionally unstable, which could possibly have negative effects on their job performance. This study also hypothesized that individuals who are open are likely to be better learners since they are curious by nature, and companies can benefit from hiring individuals of such nature as they will learn their footings within the workplace in less time compared to other individuals. Different traits may be more beneficial for candidates looking for jobs in particular fields. For example, agreeableness and extraversion traits would be beneficial for an employee working in Sales or Marketing. Cohen, Ornoy, and Keren [7] describes a study where project managers and their personality types were identified by the 16 MBTI personality types. This research identifies that there's a correlation between successful project managers and their personality types.

Finding candidates whose personality is a great fit for the role will raise employee retention rates as well since employees are likely to be happier if they are a good fit for the role not just in skill, but in personality as well. Currently, most companies assess a candidate's personality by interviewing them. This is very time consuming, and a company is only realistically able to interview a few individuals. Great candidates may fall through the cracks because of the limited time and because it is quite hard to determine an individual's personality and compatibility to the company over a short interview. This may lead to companies choosing employees that may not be the best fit for them.

A personality detection system will significantly improve the effectiveness of e-recruitment for an organization. Automating the process of assessing a candidate's personality will result in the following benefits for the recruiting organization.

- 1) Typically, companies are only able to assess a candidate's personality once they conduct interviews with the candidates. Through the personality detection system recruiters can get an assessment before the interview stage with no manual work or time involved. Overall, this system will increase the effectiveness of e-recruitment.
- 2) Recruiters are usually only able to interview a couple of candidates since there is limited time available. Hence the recruiters will have to filter candidates to a few individuals to conduct interviews. With the personality detection system recruiters can assess many more candidates without spending time to conduct one on one

interviews. Overall recruiters can make better informed choices regarding candidate selection.

### III. RELATED WORK

Majumder, Poria *et al.* [8] presents a system that detects personality via text using deep learning. In this project the existence or absence of ‘the Big Five’ is assessed using five separate binary classifiers for each trait. Each classifier is a deep CNN with document level features that’s extracted from the text itself and fed to an inner layer. The starting layers of the neural network focuses on the text in separate sentences and then eventually the text is treated as an aggregate [8]. The average accuracy achieved of each trait is as follows: Extraversion -58.09%, neuroticism – 59.38%, agreeableness – 56.71%, conscientiousness – 56.73%, openness – 62.68%.

Ren, Shen *et al.* [9] uses BERT for text semantic extraction. This research focuses on overcoming the gaps of keyword only extraction for personality detection and provides more analysis about sentiment information. This system aims to detect personality from social media posts. This system has been tested on MBTI personality types as well as the Big Five traits. The average accuracy achieved of each Big Five trait is as follows: Extraversion -79.94%, neuroticism – 80.14%, agreeableness – 80.30%, conscientiousness – 80.23%, openness – 80.35%. The test accuracy results for the MBTI personality types are E-I: 0.8175, S-N: 0.9075, T-F: 0.7931, J-P: 0.7876. These results were achieved with Bert along with CNN.

Rahman, Faisal *et al.* [10] aims to produce a comparison between various deep learning algorithms to detect personality from text. The research states that deep learning algorithms have fairly good performance however this performance may vary depending on what activation function is used. The paper concludes that the B for tanh activation function is the best among its competitors. This research also uses ‘the Big Five’ traits for testing. The F1-score of sigmoid, tanh and leaky ReLU respectively are 33.11%, 47.25%, and 49.07% [9].

Sun, Liu *et al.* [11] proposes a model named 2CLSTM that encompasses a bidirectional LSTM and along with a CNN to detect users personality based on text input. This research is tested on two different datasets, short text and long text. Sudha *et al.* [12] proposes a system ‘AdaWalk’ which is an unsupervised learning model and analyses groups of user records. This system was used to predict ‘The Big Five’ and performed well compared to other models. The results of accuracy of identifying big five characteristics range from 0.5546 – 0. 6769.

Kazameini, S. Fatehi *et al.* [13] compares bagged SVM and BERT. This research aims to develop a system that is computationally efficient, that can be used by many people. It presents a Bagged-SVM classifier that is fed with contextualized embeddings and psycholinguistic features. This project was able to achieve 58.51%.

Cui, Qi, *et al.* [14] studies different NLP approaches and evaluate them on classifying MBTI based on user’s social media posts. The project provides the following results with the relevant NLP approaches. SoftMax -17%, Naive Bayes – 26%, Regularized SVM - 33%, Deep learning – 38%.

Ontoum., Chan *et al.* [15] presents several machine

learning techniques to predict MBTI personality. This project uses CRISP-DM to guide the learning process. The best results were achieved with Recurrent neural networks with a 49.75% accuracy.

### IV. PROPOSED SOLUTION AND IMPLEMENTATION

#### A. Proposed Solution

Fig. 1 describes the feature diagram of the prototype. The user will initially enter their details as well as CV/ personal statement. This text data is preprocessed and sent to the BERT model which is finetuned on MBTI dataset. In the end the personality of that user is detected and the matching job posting details are fetched from the connected database and finally the output will be the detected personality details and the relevant job postings that correspond to the detected personality.

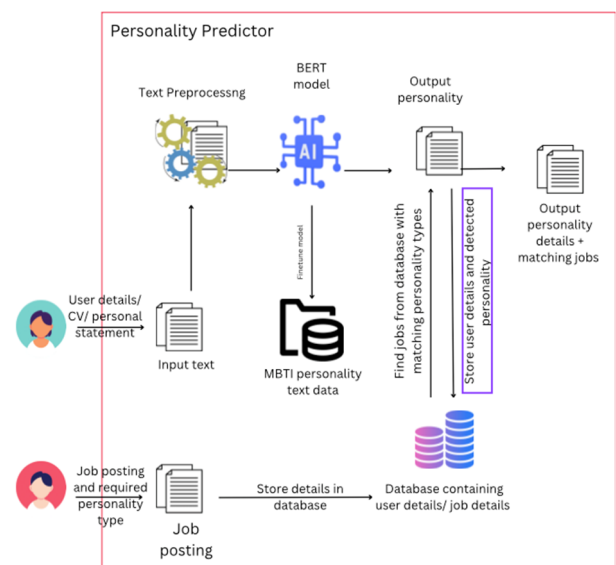


Fig. 1. Prototype feature diagram.

#### B. Dataset and Preprocessing

For this project the user uses the Kaggle MBTI personality types dataset. In order to achieve optimal results in this project, it was crucial to carefully select the dataset based on specific components, which include:

- 1) A text component, such as a comprehensive description or essay, that contains sufficient information to detect its underlying semantics. This component serves as a simulation of a CV.
- 2) The corresponding personality type inferred from the text.

This dataset contains user written content extracted from reddit; social media categorized into the MBTI personalities. This dataset was the closest the author could get to simulating an essays/ CV dataset categorized into MBTI personality types.

In this dataset, each data point represents a single post, and the features describe the contents of the post and the personality type of the user who wrote it. The tuples in this dataset consist of two elements: the contents of the post, represented as a string, and the MBTI personality type of the user who posted it, represented as a four-letter code. A Tuple will look like the following:

“I’ve always had a fascination with things like this. I think the “hard” sciences can explain a lot, but there’s still a lot of unexplained phenomena that’s worth exploring. I think one of the most interesting parts of this is trying to separate fact from fiction. There’s so much superstition and mythology surrounding this kind of thing, but I think that if we approach it with a sceptical but open mind, we might be able to make some real progress”, “INTJ”).

The dataset also includes additional information in the form of column headers. The two columns in this dataset are:

- 1) Type: The MBTI type (personality) of the user who posted the corresponding post.
- 2) Posts: The contents of the post, represented as a string, represented as a four-letter code.

Each row in the dataset corresponds to a compilation of posts, with the type column indicating the MBTI personality type of the user who posted it, and the posts column containing the contents of the post itself.

The dataset contains 8675 rows. Among these rows the 16 different personality classes are represented. However, there is an imbalance between these classes as depicted by Fig. 2. This imbalance was handled by data augmentation and pool sampling.

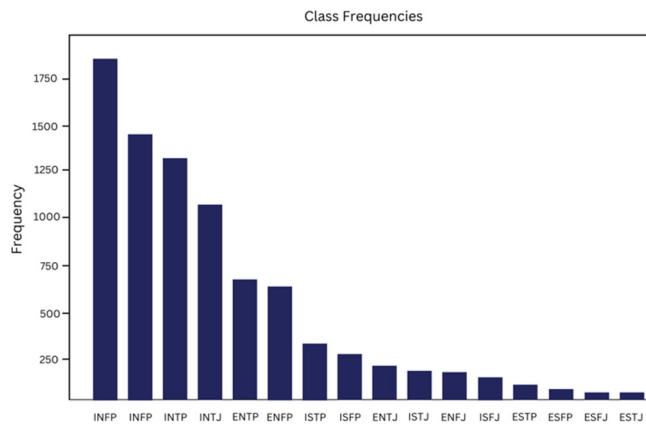


Fig. 2. Class Frequency of dataset.

Data augmentation was conducted for the classes with less frequency of occurrence, by splitting posts within a single row and recombining them with other posts of similar personality type to create a separate row of data. Posts within a single row on the dataset are separated by ‘|||’. With this data augmentation method, the author was able to combat class imbalance to a certain extent. However, with the intensity of the imbalance present in this dataset, under sampling the classes with high frequency were also necessary.

Once the class imbalance was dealt with as best possible, the dataset was preprocessed with the following steps:

- 1) Convert text to lowercase
- 2) Remove brackets, URLs, any html tags and ‘\n’ characters.
- 3) Remove any numbers.
- 4) Remove any non-ASCII characters.
- 5) Removing quotes.

### C. Model Choosing and Architecture

The author decided to use the BERT pretrained model and

conduct transfer learning for this project. BERT was chosen for this project mainly based on the following reasons:

- 1) Lack of a large dataset. The only suitable dataset for this project is relatively small and has imbalances between the sixteen classes represented within it. Once techniques such as sample pooling and data augmentation is done to reduce class imbalances the dataset can be reduced even further in size. PLMs are able to achieve high accuracy through transfer learning even through limited datasets, hence the author chose BERT PLM for this project.
- 2) BERT is pre-trained to analyze unlabeled text bidirectionally, meaning that it takes into account both the left and right contexts in all of its layers. [16]. This provides better outcomes for this project rather than a model that only processes data left to right.
- 3) BERT is freely available and easy to use.
- 4) Since the time frame for this project is limited, it’s convenient to use a PLM such as BERT to be trained well on limited time. This allows the author to train the model multiple times in multiple ways to experiment and achieve the best results.
- 5) BERT is less resource intensive when compared to its successors. Because the author was forced to train the model locally on a system that did not have an external GPU, Bert was chosen since it was still able to function with good performance.

The architecture of BERT used for this implementation is as follows. In this specific implementation, the input sequence is tokenized using the BERT tokenizer and then passed through the pre-trained BERT model to generate a sequence of hidden states. These hidden states represent the learned representations of the input sequence. Next, global average pooling is applied to the hidden state tensor and resulting pooled hidden states are then passed through a dense output layer with a sigmoid activation function.

BERT has been fine tuned to perform a multi-label classification task using the below layers:

- 1) Global Average Pooling Layer: This layer is used to reduce the dimensionality of the output hidden state tensor. The resulting pooled hidden states capture the most important features of the input sequence across all time steps.
- 2) Dense Output Layer: This layer is added to enable multi-label classification, where each output axis represents a binary classification task. The dense layer is followed by a sigmoid activation function, which produces a prediction for each output axis. The sigmoid function ensures that the output values are between 0 and 1, which can be interpreted as the probability of each label being present in the input sequence.

Overall, this model architecture is designed to use the pre-trained BERT model to learn representations of natural language inputs, and then fine-tune those representations for this multi-label classification task.

## V. RESULTS

Table I and II depict the performance of the proposed method. Table III depicts the comparison between the proposed solution and similar research.

Table 1. Accuracy of each class

Class	Accuracy
I - E	58.97%
N - S	73.64%
T - F	75.54%
J - P	61.68%
Weighted Average	72.14%

Table 2. Metrics (Weighted)

Metrics	Value
Accuracy	72.14%
Precision	81.06%
Recall	78.83%
F1 score	79.92%

Table 3. Benchmarking

Approach	Weighted Accuracy for singular class
Softmax [14]	64.21%
Naïve Bayes [14]	70.55%
SVM [14]	76.1%
LSTM [14]	78.51%
RNN [15]	83.98%
BERT finetuned (proposed)	72.14%
XGBoost [17]	97.34%

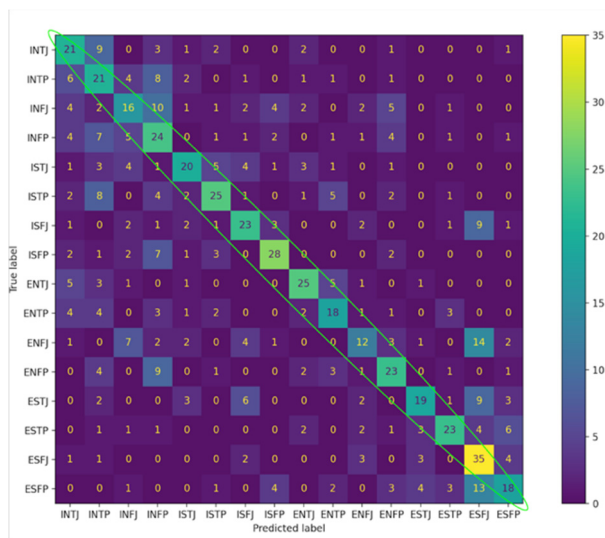


Fig. 3. Confusion matrix.

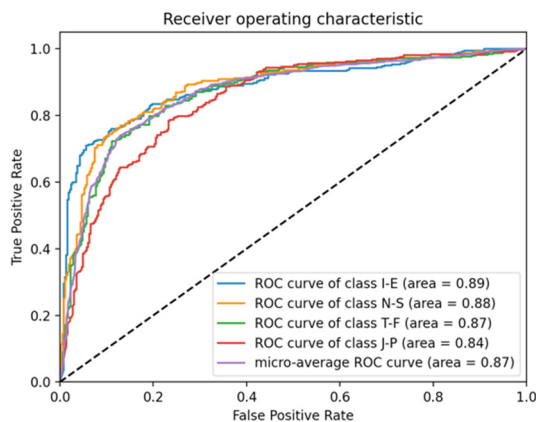


Fig. 4. ROC AUC Curve.

Fig. 3. depicts the plotted confusion matrix and Fig. 4. Depicts the ROC AUC. The AUC for these models' classes ranges from 0.72 - 0.82 which depicts that there is a relatively good separation measure within this model. We

can see that the system works reasonably well in comparison. The author believes that with future enhancements such as building a more balanced dataset and further model enhancement, performance can be increased significantly.

## VI. CONCLUSION

In this paper, the authors have provided a solution to detect a candidates MBTI personality based on their CV/personal statement for the domain of e-recruitment. This research proposes the use of PLMs, more specifically finetuned BERT for this classification task. The model was trained on the Kaggle MBTI dataset which was processed further because of present class imbalances. The model was able to achieve a weighted average of 72.14% and performs relatively well when compared with other models that were benchmarked on the same dataset. We believe that this avenue of research is very relevant and could provide to be very useful in the field of e-recruitment. Our next objective is to research and develop this model further and perform further comparisons and adjustments.

## CONFLICT OF INTEREST

The authors declare no conflict of interest.

## AUTHOR CONTRIBUTIONS

Himaya Perera and Lakshan Costa conducted the research; Both authors analyzed the data; Himaya Perera wrote the paper. All authors had approved the final version.

## ACKNOWLEDGMENT

The authors of this paper firstly would like to thank the Informatics Institute of Technology Sri Lanka for the guidance, support and resources given for the completion of this research. The authors would also like to thank each individual that helped in the process of concluding this research.

## REFERENCES

- [1] L. Barber. (2006). *e-Recruitment Developments*. Institute for Employment Studies, University of Sussex, 5, [Online]. Available: <https://www.employment-studies.co.uk/system/files/resources/files/mp63.pdf>
- [2] K. Dumper, W. Jenkins, A. Lacombe, M. Lovett, and M. Perimutter, "10.1 What is Personality?" *Introductory Psychology*, p. 481, 2019.
- [3] D. H. Pelt, D. van der Linden, C. S. Dunkel, and M. P. Born, "The general factor of personality and job performance: Revisiting previous meta-analyses," *International Journal of Selection and Assessment*, vol. 25, no. 4, pp. 333-346, 2017.
- [4] F. Schmidt. (2016). The validity and utility of selection methods in personnel psychology: practical and theoretical implications of 100 years of research findings" 2, [Online]. Available: <https://www.researchgate.net/publication/309203898>
- [5] K. Cherry, "What are the big 5 personality traits?" *Verywell Mind*, March 2023.
- [6] M. Barrick and M. Mount, "The big five personality dimensions and job performance: A meta-analysis," *Personnel Psychology*, vol. 44, issue 1, pp.1-26, March 1991.
- [7] Y. Cohen, H. Ornoy, and B. keren, "MBTI personality types of project managers and their success: A field survey," *Project Management Journal*, vol.44, issue 3, pp. 78-87, June 2013.
- [8] N. Majumder, S. Poria, A. Gelbukh, and E. Cambria, "Deep learning-based document modeling for personality detection from text," *IEEE Intelligent Systems*, vol. 32, no. 2, pp. 74-79, 2017.
- [9] Z. Ren, Q. Shen, X. Diao, and H. Xu, "A sentiment-aware deep learning approach for personality detection from text," *Information Processing & Management*, vol. 58, issue 3, 2021.

- [10] M. A. Rahman, A. Al Faisal, T. Khanam, M. Amjad, and M. S. Siddik, "Personality detection from text using convolutional neural network," presented at 2019 1st International Conference on Advances in Science, Engineering and Robotics Technology (ICASERT), Dhaka, Bangladesh, 3rd–5th May, 2019.
- [11] X. Sun, B. Liu, J. Cao, J. Luo, and X. Shen, "Who am I? Personality detection based on Deep Learning for texts," presented at 2018 IEEE International Conference on Communications (ICC), Kansas City, Missouri, USA, 20th-24th May, 2018.
- [12] G. Sudha, S. K. K, S. J. S, N. D, S. S, and K. T. G, "Personality prediction through CV analysis using machine learning algorithms for automated E-recruitment process," presented at 2021 4th International Conference on Computing and Communications Technologies (ICCCT), Chennai, India, 16th–17th December, 2021.
- [13] A. Kazameini, S. Fatehi, Y. Mehta, S. Eetemadi, and E. Cambria, "Personality trait detection using bagged SVM over Bert Word embedding ensembles," arXiv.org, October 2020, [Online]. Available: <https://arxiv.org/abs/2010.01309>.
- [14] B. Cui and C. Qi. (2017). Survey analysis of machine learning methods for natural language," Stanford University. [Online]. Available: <http://cs229.stanford.edu/proj2017/final-reports/5242471.pdf>
- [15] S. Ontoum and J. H. Chan, "Personality type based on Myers-Briggs Type Indicator with text posting style by using traditional and deep learning," arXiv.org, January 2022. [Online]. Available: <https://arxiv.org/abs/2201.08717>
- [16] J. Duan, H. Zhao, Q. Zhou, M. Qiu, and M. Liu, "A study of pre-trained language models in Natural Language Processing," presented at 2020 IEEE International Conference on Smart Cloud (SmartCloud), Washington, District of Columbia, USA, 6th-8th November, 2020.
- [17] A. S. Khan, H. Ahmad, M. Zubair, F. Khan, A. Arif, and H. Ali, "Personality classification from online text using machine learning approach," *International Journal of Advanced Computer Science and Applications*, vol. 11, no. 3, 2020.

Copyright © 2024 by the authors. This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited ([CC BY 4.0](https://creativecommons.org/licenses/by/4.0/)).