

# **AUTOMATIC PERSONALITY DETECTION USING DEEP LEARNING MODELS IN USER-GENERATED CONTENT**

A Thesis Submitted  
In Partial Fulfilment of the Requirements for the  
Degree of

**DOCTOR OF PHILOSOPHY**

by

**DIPIKA JAIN**  
(2k20/PHDCO/505)

Under the supervision of

**Dr. AKSHI KUMAR**  
Goldsmiths, University of London, United Kingdom  
and  
**Dr. ROHIT BENIWAL**  
Delhi Technological University, New Delhi, India



**Department of Computer Science and Engineering**

**DELHI TECHNOLOGICAL UNIVERSITY**  
**(Formerly Delhi College of Engineering)**  
**Shahbad Daulatpur, Main Bawana Road, Delhi-110042. India**

**August, 2024**

## ACKNOWLEDGEMENT

I take the opportunity to express my sincere gratitude to **Dr. Akshi Kumar**, Associate Professor and Director PGR, Department of Computing, Goldsmiths, University of London, London, United Kingdom and **Dr. Rohit Beniwal**, Assistant Professor, Department of Computer Science and Engineering, Delhi Technological University, Delhi, for providing valuable guidance, consistent support, constant encouragement, motivation and direction throughout the work.

I am very grateful to **Prof. Rahul Katarya**, DRC Chairman, Department of Computer Science and Engineering, Delhi Technological University, Delhi, for his support, guidance, and encouragement.

I am very thankful to **Prof. Vinod Kumar**, Head of the Department, Computer Science and Engineering, Delhi Technological University, Delhi, for his continued support and encouragement.

I am also thankful to all other faculty members of the Dept. of Computer Science and Engineering, Delhi Technological University, Delhi, for the motivation and inspiration. I am also thankful to all non-teaching staff at DTU, and who have helped me directly or indirectly in accomplishing this research plan.

**Dipika Jain**

2K20/PHD/CO/505

Teaching and Research Scholar

Department of Computer Science and Engineering

E-mail: [dipikajain\\_2k20phdco505@dtu.ac.in](mailto:dipikajain_2k20phdco505@dtu.ac.in)

## CANDIDATE DECLARATION

I, Dipika Jain, full-time research scholar (2k20/PHDCO/505) in the Department of CSE, DTU, hereby declared that the thesis titled *Automatic Personality Detection using Deep Learning models for user-generated content*, which is being submitted towards the fulfilment of the requirements for the degree of Doctor of Philosophy in Computer Science & Engineering of Delhi Technological University, New Delhi is a record of bonafide research work carried out by me. I further declare that this work is based on original research and has not been submitted to any university or institution for any degree or diploma.

**Dipika Jain**

2k20/PHDCO/505

Teaching and Research Scholar

Dept. of Computer Science and Engineering

Delhi Technological University, New Delhi, India

**E-mail:** dipikajain\_2k20phdco505@dtu.ac.in



**DELHI TECHNOLOGICAL UNIVERSITY**  
Shahbad Daulatpur, Main Bawana Road, Delhi 110042.

## **SUPERVISOR(S) CERTIFICATE**

Date: 23-08-2024

This is to certify that the work embodied in this thesis entitled *Automatic personality detection using deep learning models in user-generated content*” done by Dipika Jain, Roll no. 2K20/PHD/CO/505 as a full-time scholar in the Department of Computer Science and Engineering, Delhi Technological University is an authentic work carried out by her under our guidance.

This work is based on original research and the matter embodied in this thesis has not been submitted earlier for the award of any degree or diploma to the best of my knowledge and belief.

### **Supervisors**

**Dr. Rohit Beniwal**

Assistant Professor,  
Dept. of Computer Science & Engineering,  
Delhi Technological University  
New Delhi, India  
[rohit.beniwal@dtu.ac.in](mailto:rohit.beniwal@dtu.ac.in)

**Dr. Akshi Kumar**

Associate Professor and Director PGR,  
Department of Computing,  
Goldsmiths, University of London  
London, United Kingdom  
[akshi.kumar@gold.ac.uk](mailto:akshi.kumar@gold.ac.uk)

## ABSTRACT

This thesis titled *"Automatic Personality Detection Using Deep Learning Models in User-Generated Content"* presents a cutting-edge exploration into the realm of artificial intelligence, focusing on the application of deep learning models to discern and analyze personality traits from vast arrays of user-generated content (UGC). This work stands at the intersection of computer science, psychology, and data analytics, endeavouring to bridge the gap between quantitative AI techniques and qualitative psychological insights. Central to the research is the utilization of advanced machine learning models, particularly transformer-based architectures like BERT and GPT, renowned for their ability to process and understand natural language at a granular level. These models are applied to a variety of data sources, including social media posts, blog entries, and online interactions, to extract and classify personality indicators according to well-established frameworks such as the Big Five personality traits and the Myers-Briggs Type Indicator (MBTI).

The research is innovative in its approach, using not just textual analysis but also exploring multimodal data, including video and audio, to capture the nuances of personality expression that text alone might miss. By integrating deep learning with psychological theories, the study aims to develop models that not only predict personality types with high accuracy but also understand the underlying human behaviours and emotional states that drive online interactions.

One of the significant contributions of this research is its focus on cross-linguistic and cross-cultural applicability. Recognizing the diversity of global online communities, the study tests and validates models across different languages and cultural contexts, enhancing the generalizability and usefulness of the findings. This aspect is crucial, considering the potential for AI systems to perpetuate biases if not properly trained on diverse datasets. The thesis also explores the practical application of its findings in real-life domain, demonstrating how personality detection can be utilized to improve recruitment processes by aligning candidates' personalities with organizational culture. In detailing the results, the research discusses the performance of various models tested, providing comparative analyses that highlight the strengths and limitations of each approach.

In conclusion, this thesis not only advances the field of AI and personality psychology but also contributes to the broader discourse on the integration of technology and human-centric sciences. By forging a link between computational models and psychological theories, it lays a foundation for future innovations that respect human diversity and foster more personalized and meaningful interactions in the digital age.

## LIST OF PUBLICATIONS

### JOURNAL:

1. Kumar, A., Beniwal, R, Jain, D. (2023). Personality Detection using Kernel-based Ensemble Model for leveraging Social Psychology in Online Networks *ACM Transactions on Asian and Low-Resource Language Information Processing (ACM TALLIP)*- <https://doi.org/10.1145/3571584> [SCIE-Impact Factor: 1.472] PUBLISHED
2. Kumar, A., Jain, D., & Beniwal, R. (2023). HindiPersonalityNet: Personality Detection in Hindi Conversational Data using Deep Learning with Static Embedding. *ACM Transactions on Asian and Low-Resource Language Information Processing*. <https://doi.org/10.1145/3625228>[SCIE-Impact Factor: 1.472] PUBLISHED
3. Jain, D., & Kumar, A., (2024). AI unveiled personalities: Profiling optimistic and pessimistic attitudes in Hindi dataset using transformer-based models, *Expert Systems*, John Wiley & Sons <https://doi.org/10.1111/exsy.13572> [SCIE-Impact Factor: 2.812], PUBLISHED
4. Jain, D., Beniwal, R. & Kumar, A. (2024). Advancements in Personality Detection: Unleashing the Power of Transformer-Based Models and Deep Learning with Static Embeddings on English Personality Quotes. *International Journal of All Research Education & Scientific Methods, (IJARES M)*, <https://doi.org/10.56025/IJARESM.2023.1201242235> [UGC-Care] PUBLISHED
5. Jain, D., & Kumar, A., Hyper-Personalized Employment in Urban Hubs: Multimodal Fusion Architectures for Personality-Based Job Matching. *Neural Computing and Applications*, Springer. (Revisions Submitted)
6. Kumar, A. & Jain, D., Enhancing Decision-Making Across Domains with MBTI and Transformer Models: Applications from General AI to Health Informatics, *Scalable Computing: Practice and Experience* (Communicated)

7. Jain, D., & Kumar, A., *Machine Learning for Predicting Parenting Styles through MBTI Profiles in Hindi-Speaking Populations* (**Communicated**)
8. Kumar, A. & Jain, D., *EmoMBTI-Net: Introducing and Leveraging a Novel Emoji Dataset for Personality Profiling with Large Language Models* (**Communicated**)

#### **CONFERENCE:**

1. Jain, D., Kumar, A., Beniwal, R. (2021). Personality BERT: A Transformer-Based Model for Personality Detection from Textual Data. *In Proceedings of International Conference on Computing and Communication Networks. Lecture Notes in Networks and Systems*, vol 394. Springer, Singapore. [https://doi.org/10.1007/978-981-19-0604-6\\_48](https://doi.org/10.1007/978-981-19-0604-6_48)
2. Jain, D., Beniwal, R., Kumar, A. (2023). ByaktitbaNet: Deep Neural Network for Personality Detection in Bengali Conversational Data. *In Proceedings of Fourth Doctoral Symposium on Computational Intelligence. DoSCI 2023. Lecture Notes in Networks and Systems*, vol 726. Springer, Singapore. [https://doi.org/10.1007/978-981-99-3716-5\\_57](https://doi.org/10.1007/978-981-99-3716-5_57) [**Best Paper Awarded**]

## TABLE OF CONTENTS

<b>1. Introduction .....</b>	<b>18-26</b>
1.1. Challenges in Detecting Personality Automatically .....	19
1.2. Advances in Personality Detection: The Role of Transformer and Large Language Models (LLMs) .....	20
1.3. Personality Models and Traits Overview .....	21
1.3.1. The Big Five Personality Traits (OCEAN) .....	21
1.3.2. Myers-Briggs Type Indicator (MBTI) .....	22
1.3.3. Extraversion, Introversion, and Ambiversion .....	22
1.4. Statement of Research and Research Objectives .....	23
1.5. Primary Contribution of the Research .....	24
1.5.1. Contribution to the Knowledge Base .....	24
1.6. Organization of the Thesis .....	24
1.7. Chapter Summary .....	26
 <b>2. Literature Review .....</b>	 <b>27-36</b>
2.1. Traditional Methods of Personality Detection .....	27
2.2. Advancements in Personality Detection through Deep Learning and Data Analytics .....	27
2.3. Major findings from Literature Review .....	34
2.4. Chapter Summary .....	35
 <b>3. Research Datasets .....</b>	 <b>37- 58</b>
3.1. Significance of the Diverse Spectrum of Datasets .....	38
3.2. English Language Data .....	39
3.2.1. First Impression Chalearn (English, Public) .....	40
3.2.2. Kaggle_MBTI (English, Public) .....	41
3.2.3. Personality_Quotes (English, Curated) .....	42
3.2.4. DM_MBTI (English, Reannotated) .....	43



3.2.5.	Emo_MBTI (English, Reannotated) .....	45
3.3.	Hindi Language Data .....	47
3.3.1.	Shaksiyat (Hindi, Curated) .....	48
3.3.2.	Vishesh_Charitr (Hindi, Curated) .....	49
3.3.3.	Manobhav (Hindi, Reannotated) .....	50
3.3.4.	Parvarish (Hindi, Curated) .....	51
3.4.	Bangla Language Data: Byaktitba (ব্যক্তিত্ব) .....	52
3.5.	Validation and Testing of Datasets .....	53
3.6.	Preprocessing of Datasets .....	54
3.6.1.	Text Cleaning and Normalization .....	54
3.6.2.	Handling of Emojis .....	55
3.6.3.	Audio and Video Preprocessing .....	56
3.6.4.	Language-Specific Preprocessing for Hindi and Bangla .....	57
3.7.	Chapter Summary .....	58

#### **4. RO1: Developing and Validating Cross-Linguistic Personality Detection**

##### **Models Using Native Language Datasets and Emoji Integration ..... 59-72**

4.1.	Classical Deep Learning Models .....	60
4.2.	Word Embeddings .....	61
4.3.	Transformer-Based models .....	63
4.4.	PersonalityBERT (English) .....	65
4.5.	HindiPersonalityNet (Hindi) .....	65
4.6.	ByaktitbaNet (Bangla) .....	67
4.7.	Transformer-based Models for Personality_Quotes .....	69
4.8.	Chapter Summary .....	71

#### **5. RO2: Integrating Psychological Theories into Deep Learning Frameworks for Multi-Language and Cultural Contexts ..... 73-80**

5.1. Kernel-based Soft Voting Ensemble Model for Personality Detection (KBSVE-P).....	73
5.1.1. Support Vector Machine (SVM) .....	74
5.1.2. Voting .....	75
5.2. Transformer-Based Models for Decision-Maker-MBTI Type .....	76
5.3. Machine Learning for Predicting Parenting Styles via MBTI .....	78
5.3.1. Parenting Styles and MBTI Profiles .....	78
5.4. Chapter Summary .....	79
 <b>6. RO3: Profiling Emotional Dispositions and Evaluating Real-World Applications .....</b>	<b>81-96</b>
6.1. Transformer-Based Models for Emotional Profiling .....	81
6.1.1. Overall Human Psychological Make-up .....	82
6.1.2. Human Psychology Chakra System .....	83
6.1.3. Proposed Stacked mDeBERTa Model .....	87
6.2. Fusion Model: XceptionResNet with BERT for Chalearn (Multimodal) Dataset .....	88
6.2.1. Deep Neural Architectures .....	89
6.2.2. Multimodal Fusion Architectures .....	89
6.3. Emo-MBTI: Mapping Emojis to MBTI Personalities .....	91
6.3.1. EmoMBTI-Net: The Emoji-Based MBTI Personality Prediction Model .....	92
6.4. Chapter Summary .....	95
 <b>7. Results and Discussion .....</b>	<b>97-134</b>
7.1. Performance Metrics used in Research .....	97
7.1.1. Model Training and Validation .....	98
7.2. Results of Research Objective 1 .....	99
7.2.1. Model: PersonalityBERT .....	99

7.2.2. Model: HindiPersonalityNet .....	99
7.2.3. Model: ByaktitbaNet .....	103
7.2.4. Transformer-based Models for Personality_Quotes .....	104
7.3. Results for Research Objective 2 .....	109
7.3.1. Model: KBSVE-P .....	110
7.3.2. DM_MBTI .....	117
7.3.3. ParvarishNET .....	119
7.4. Results for Research Objective 3 .....	121
7.4.1. Model: Transformer-based for Manobhav Dataset .....	121
7.4.2. Multimodal fusion architecture .....	124
7.4.3. EmoMBTI-NET .....	129
7.5. Limitations of the study .....	132
7.5.1. Data Bias and Generalizability .....	133
7.5.2. Interpretability Challenges .....	133
7.5.3. Data Quality and Annotation Consistency .....	133
7.5.4. Scalability and Computational Resources .....	133
7.5.5. Ethical and Privacy Concerns .....	134
7.6. Chapter Summary .....	134
<b>8. Conclusion and Future Work .....</b>	<b>135-141</b>
8.1. Summary of Key Contributions .....	135
8.1.1. Results of Research Objective-1 .....	136
8.1.2. Results of Research Objective-2 .....	137
8.1.3. Results of Research Objective-3 .....	138
8.2. Social Impact .....	139
8.3. Recommendations for Future Research .....	140
8.4. Closing Thoughts .....	141
<b>References .....</b>	<b>142-151</b>

<b>Appendix A:</b> List of publications with Proofs .....	<b>152-155</b>
<b>Appendix B:</b> Plagiarism Report .....	<b>156</b>
<b>Appendix C:</b> Curriculum Vitae .....	<b>157-159</b>

## LIST OF TABLES

**Table 2.1:** Literature related to Personality Detection in English Datasets  
**Table 2.2:** Literature on Personality Detection in Hindi Dataset  
**Table 2.3:** Literature on personality detection in other non-English Languages

**Table 3.1:** An overview of Datasets  
**Table 3.2:** *Personality\_Quotes* Dataset statistics  
**Table 3.3:** Shakhsiyat (शख्सियत) Dataset statistics  
**Table 3.4:** Snapshot of Shakhsiyat (शख्सियत) dataset  
**Table 3.4:** Byaktitba (ব্যক্তিত্ব) Dataset statistics

**Table 7.1.** Comparative Performance of Hindi Personality Detection Models  
**Table 7.2.** Comparative Analysis of Hindi and Bangla Personality Detection  
**Table 7.3.** Comparative Analysis of various deep learning models with BioWordVec  
**Table 7.4.** Accuracy comparison of embedding & deep learning combinations  
**Table 7.5.** Performance of ByaktitbaNet Model  
**Table 7.6.** Performance comparison with another dataset  
**Table 7.7.** Accuracies for classical deep learning models with various static embeddings  
**Table 7.8.** Performance of transformer-based models on *personality\_quotes* dataset  
**Table 7.9.** Performance of Transformer-Based Models on *Personality\_Quotes* Dataset  
**Table 7.10.** Performance of different SVM kernels: Kaggle\_MBTI  
**Table 7.11.** Performance of ensemble of SVM kernels aggregated with different voting techniques  
**Table 7.12.** Comparison with the State-of-Art (Accuracy): Kaggle\_MBTI  
**Table 7.13.** KBSVE-P on विशेषचरित्र\_\_MBTI dataset  
**Table 7.14.** Performance of all Transformer-based Models on DM-MBTI dataset  
**Table 7.15.:** Performance of ParvarishNET  
**Table 7.16:** Comparison of the results of parenting styles  
**Table 7.17:** Performance of different models over Parvarish Dataset  
**Table 7.18.:** Performance of all Transformer-based Models on मनोभाव dataset  
**Table 7.19.** Model Hyperparameters  
**Table 7.20.** Performance across modalities  
**Table 7.21.** Performance of Text Modality Models  
**Table 7.22.** Performance of Audio-generated Images Models  
**Table 7.23.** Performance of Video frames  
**Table 7.24.** Performance of multimodal decision models  
**Table 7.25.:** Emoji mapping Rouge scores using LLMs  
**Table 7.26.:** Accuracy of different models  
**Table 7.27.:** F1- scores of different models  
**Table 7.28.:** Hyperparameters used

## LIST OF FIGURES

- Fig. 3.1.** Transcripts of the dataset  
**Fig.3.2.** Frames of the Dataset  
**Fig. 3.3.** Personalities distribution in the MBTI dataset  
**Fig.3.4.** *Personality\_Quotes* Dataset snapshot  
**Fig.3.5.** Decision maker-MBTI Personality Trait Mapping  
**Fig.3.6.** Snapshot of DM-MBTI Dataset  
**Fig.3.7.** MBTI Posts relabelled with Emojis  
**Fig.3.8.** Snapshot of Emoji Labelled MBTI Dataset  
**Fig.3.9.** Sample of Hindi dataset विशेष चरित्र\_MBTI  
**Fig.3.10** ‘मनोभाव’ Dataset snapshot  
**Fig. 3.11.** Mapping MBTI personality Traits to Pessimistic-Optimistic Attitudes  
**Fig.3.12.** Snapshot of Parvarish dataset  
**Fig.3.13** Snapshot of Byaktitba (ব্যক্তিত্ব) Dataset  
**Fig. 4.1.** Architecture of PersonalityBERT model  
**Fig.4.2.** Deep learning and embedding employed  
**Fig.4.3.** The HindiPersonalityNet Model  
**Fig.4.4.** The personality continuum spectrum  
**Fig.4.5.** ByaktitbaNet Model  
**Fig.4.6.** Word Embeddings with Deep Learning Models  
**Fig.4.7.** Personality Prediction using Transformer-based Models  
**Fig. 5.1.** Kernel-based Ensemble Model (KBSVE-P) evaluated on English MBTI and Hindi विशेष चरित्र\_MBTI datasets  
**Fig.5.2.** Voting Classifiers (Soft Voting, Hard Voting, and Weighted-Hard Voting)  
**Fig.5.3.** MBTI Decision-Maker Type Recognition using Transformer Models evaluated on re-annotated English MBTI dataset  
**Fig.5.4.** Workflow of different Models for Personality Prediction  
**Fig.6.1.** The Human Psychology Chakra System  
**Fig.6.2.** Emotional Attitude prediction using transformer-based models  
**Fig.6.3.** XceptionResNet + BERT Structure  
**Fig.6.4.** The MBTIEmoNet Model  
**Fig.7.1.** Accuracy Comparison with State-of-the-art  
**Fig.7.2.** Performance comparison of Indian Language personality detection  
**Fig.7.3.** F1-score comparison  
**Fig.7.4.** Model accuracy curve for 100 epochs  
**Fig.7.5.** Comparison of f1 scores of transformer-based models  
**Fig.7.6.** Confusion matrix for the top-performing model, ELECTRA  
**Fig.7.7.** Comparison of accuracy of all models  
**Fig.7.8.** Comparison of F1 Scores of Transformer-Based Models  
**Fig. 7.9.** Accuracies of different SVM kernels: Kaggle\_MBTI  
**Fig. 7.10.** F1-Score of different SVM kernels: Kaggle\_MBTI

**Fig. 7.11.** Accuracies of voting techniques: Kaggle\_MBTI  
**Fig. 7.12.** F1-Score of different Voting techniques: Kaggle\_MBTI  
**Fig. 7.13.** Comparison with the existing work in Kaggle\_MBTI  
**Fig.7.14.** F1 comparison of Transformer-based Models  
**Fig.7.15.** Confusion matrix for the best performing Transformer model, DeBERTa  
**Fig.7.16.:** Confusion matrix of the ParvarishNET  
**Fig.7.17.:** F1 comparison of Transformer-based Models  
**Fig.7.18.** ROC curve for Stacked mDeBERTa  
**Fig.7.19.** ROC curve for IndicBERT  
**Fig.7.20.** ROC curve for mDeBERTa  
**Fig.7.21.** ROC curve for XLM-RoBERTa  
**Fig.7.22.** ROC curve for mBERT  
**Fig. 7.23.** F1 score comparison of transformer-based models  
**Fig.7.24.** Model Comparison for Personality Analysis: Accuracy of Multimodal Data Processing

## LIST OF ABBREVIATIONS

Abbreviations	Full forms
AI	Artificial Intelligence
UGC	User-Generated Content
NLP	Natural Language Processing
BERT	Bidirectional Encoder Representations from Transformers
GPT	Generative Pre-Trained Transformers
LLM	Large Language Models
MBTI	Myers Briggs Type Indicator
I/E	Introversion – Extroversion
N/S	Intuition – Sensing
T/F	Thinking – Feeling
J/P	Judging – Perceiving
CLS	Start of a Sequence
SEP	End of a Sequence
PAD	Padding
CNN	Convolutional Neural Network
GRU	Gated Recurrent Unit
LSTM	Long Short-term Memory
BiLSTM	Bidirectional Long Short-term Memory
GloVe	Global Vectors
DeBERTa	Decoding- enhanced BERT with Disentangled Attention
ELECTRA	Efficiently Learning an Encoder that classifies token Replacements Accurately
MLM	Masked Language Model
Meta OPT	Metas Open Pre-Trained Transformer Language Models
RNN	Recurrent Neural Network
KBSVE-P	Kernel based soft voting ensemble model for personality detection
SVM	Support Vector Machine
RBF	Radial Basis Function
TF-IDF	Term Frequency – Inverse Document Frequency
mBERT	Multilingual BERT
mDeBERTa	Multilingual DeBERTa
mRoBERTa	Multilingual RoBERTa
RoBERTa	Robustly Optimized BERT approach
NLI	Natural Language Inference
ROC-AUC	Area under Receiver Operating Characteristic Curve
MCC	Matthews Correlation Coefficient
MSE	Mean Squared Error



RMSE	Root Mean Squared Error
SOTA	State-of-the-art

# CHAPTER 1

## INTRODUCTION

Personality is a fundamental psychological concept that encompasses the enduring patterns of thoughts, emotions, and behaviours that characterize an individual's unique way of interacting with the world. It reflects a complex interplay of biological, psychological, and social factors [1], influencing how individuals perceive the world and themselves, interact with others, and cope with life's challenges. The study of personality is crucial across various disciplines, offering insights into human behaviour that are essential for psychological health, educational success, and effective interpersonal interactions.

The relevance of personality studies extends into the burgeoning field of artificial intelligence (AI), particularly in areas involving human-computer interaction and personalized AI services. AI systems that can predict or adapt to users' personalities have applications in enhancing user engagement, customizing learning experiences, improving mental health interventions, and refining marketing strategies. By incorporating personality assessment, AI technologies can better serve diverse user needs, providing more tailored and intuitive user experiences. Despite the extensive application of personality theory in various fields, traditional methods of personality assessment, such as questionnaires and observational studies, have limitations, particularly in terms of scalability, real-time analysis, and coverage of diverse populations. These methods often rely on self-reported data, which can be biased, and they do not easily accommodate the continuous monitoring or analysis of large-scale behavioural data.

The digital age has brought a proliferation of user-generated content (UGC), such as social media posts, blogs, and video recordings, which provide a rich, unstructured, and real-time dataset that reflects the diverse personalities of global users. This content offers a unique opportunity to observe the natural expression of thoughts, feelings, and behaviours across varied contexts and cultures. However, effectively harnessing this vast amount of data to analyze personality poses significant challenges that require advanced computational techniques.

Traditional analytics tools are often inadequate for processing the volume and complexity of UGC. Deep learning [2] and other machine learning techniques have the potential to transform this unstructured data into valuable insights about human personality. These models can learn from large datasets to identify patterns and nuances in language use, emotional expression, and online behaviour that correlate with known personality traits.

However, current research in applying these advanced computational techniques to personality analysis is still in its early stages. There is a critical need to develop more sophisticated models that can accurately interpret the subtleties of human personality from digital footprints. This involves not only improving the accuracy and efficiency of these models but also ensuring they are capable of ethical

considerations such as privacy, consent, and bias mitigation. Recent advancements in AI, particularly in the field of natural language processing (NLP)[3,4], have introduced models like transformers and embedding techniques that are well-suited for analyzing the subtleties of human language and behaviour encoded in UGC. Transformers, which include models like BERT (Bidirectional Encoder Representations from Transformers) [5] and GPT (Generative Pre-trained Transformer) [6], are particularly effective due to their ability to understand context within text, making them ideal for personality detection tasks where nuances in language can indicate different traits. These models are pre-trained on vast corpora of text and then fine-tuned for specific tasks such as sentiment analysis [7-11], which can be analogous to detecting emotional dispositions in personality. By leveraging embeddings, which convert text into numerical data that machines can understand, transformers can capture and analyze patterns in language usage that reflect underlying personality traits.

Moreover, most existing studies and models focus predominantly on data from English-speaking populations or high-resource environments. This limitation underscores the importance of developing cross-linguistic and cross-cultural models that can provide more inclusive and representative insights. Such models would allow for a more comprehensive understanding of global personality patterns, contributing to psychology, education, and AI in more culturally diverse ways.

In conclusion, while the study of personality has broad applications and significant importance, the field stands at a pivotal juncture where the integration of advanced computational methods with traditional psychological theories could vastly expand its impact and utility. By bridging this research gap, scientists and practitioners can unlock a deeper, more nuanced understanding of human personality that leverages the full potential of the digital era's data-rich landscape.

### **1.1. Challenges in Detecting Personality Automatically**

Detecting personality automatically presents a triad of challenges intricately connected to the multifaceted nature of personality, the heterogeneity of user-generated content (UGC), and the technological hurdles associated with data processing.

The complexity of human personality lies in its inherent variability among individuals, which defies one-size-fits-all analytical approaches. Each person's unique set of traits makes uniform analysis an ambitious task. This complexity is compounded by context sensitivity, as individuals may exhibit different facets of their personality in varying situations, leading to a dynamic that static assessments cannot easily capture. For example, someone may appear introverted in professional settings but display extroversion at social gatherings. Moreover, the subtleties of personality, such as an individual's proclivity for humour or propensity for empathy, often manifest in subtle or subconscious expressions that require sophisticated detection methods to decode accurately.

The variability in UGC adds another layer of complexity. UGC spans multiple formats, including text posts, videos, and audio recordings each calling for

distinct analytical approaches. Textual analysis may focus on linguistic cues through natural language processing, whereas video and audio analysis might involve sentiment analysis and voice modulation recognition to infer personality traits. However, these various formats suffer from quality inconsistencies. A high-quality blog post might offer rich insights into an individual's openness, while a low-resolution video may obscure valuable non-verbal cues, leading to data interpretation challenges. Additionally, cultural and linguistic differences in UGC necessitate the creation of localized models to ensure accuracy. The emotional resonance of certain words or expressions might vary dramatically across cultures, thus requiring models that are finely tuned to the cultural contexts they analyze.

Lastly, the technical challenges are formidable. The sheer data volume, with an exponential increase in UGC, creates scalability issues for processing and analysis. Traditional algorithms that may work well with structured, smaller datasets struggle to cope with the vast and unstructured nature of UGC, where a fixed format is often lacking. For instance, a model trained on structured interview responses may falter when applied to the colloquial and diverse structures of social media communication. Integration difficulties also arise as combining data from disparate sources and formats demands sophisticated solutions. A unified view of personality might require integrating insights from a user's text posts with their online shopping behaviour and physical activity tracker data, each format requiring different processing techniques to be meaningfully synthesized.

In essence, the quest to automatically detect personality in the digital realm is as challenging as it is intriguing, necessitating advanced solutions that are as nuanced and adaptable as the human personality itself.

## **1.2. Advances in Personality Detection: The Role of Transformer and Large Language Models (LLMs)**

Addressing the multifaceted challenges of automatically detecting personality traits within user-generated content (UGC), Large Language Models (LLMs)[12][13] such as Meta's OPT[14] and Google's BERT[15] represent the cutting edge of AI innovation. These models, built upon transformer architectures, employ intricate attention mechanisms that precisely parse through complex linguistic patterns, teasing out subtle cues indicative of various personality traits. Their extensive pre-training on vast, diverse language data grants them an acute sensitivity to context, crucial for interpreting the fluid expressions of personality that vary across different scenarios in human communication.

The research presented uses advanced AI technology, specifically deep learning models, to better understand personality traits from things people post online in different languages. It combines insights from psychology with sophisticated computer algorithms to analyze how people express their emotions and behaviours. The goal is to see how well these AI models can identify personality traits across various settings like social media platforms, schools, or work places and in multiple languages. This could help tailor experiences and interactions more effectively in these environments based on people's individual personalities.

### 1.3. Personality Models and Traits Overview

This thesis explores the complex and multifaceted realm of human personality by employing a blend of multiple well-established personality models. Each model offers a unique lens through which personality traits can be analyzed and understood, thereby contributing to a more holistic and comprehensive exploration of individual differences. The foundational models that guide this research include the Myers-Briggs Type Indicator (MBTI) [16-28], the Big Five Personality Traits (commonly known as OCEAN) [29-40], and the distinct behavioural tendencies of Extraversion, Introversion, and Ambiversion. Together, these models provide a robust framework for examining the multifaceted nature of personality.

The Big Five model captures the broad dimensions that define an individual's consistent patterns of thoughts, emotions, and behaviours, while the MBTI offers a more nuanced categorization of personality types based on specific preferences in how people perceive the world and make decisions. Additionally, the focus on Extraversion, Introversion, and Ambiversion allows for a deeper understanding of how individuals navigate social environments, providing insights into the dynamic interplay between solitary and social tendencies.

By using these diverse models, this research not only aims to identify and analyze distinct personality traits but also to explore the practical implications of these traits in various real-world contexts.

#### 1.3.1. The Big Five Personality Traits (OCEAN)

The Big Five model, also known as OCEAN, is one of the most widely accepted frameworks for understanding personality. It encompasses five broad dimensions of personality:

- **Openness to Experience:** Reflects an individual's creativity, curiosity, and willingness to embrace new ideas and experiences.
- **Conscientiousness:** Indicates a person's level of organization, dependability, and goal-oriented behaviour.
- **Extraversion:** Measures sociability, assertiveness, and the tendency to seek stimulation in social interactions.
- **Agreeableness:** Represents the degree of warmth, kindness, and cooperativeness an individual exhibits.
- **Neuroticism:** Describes emotional stability and the propensity to experience negative emotions like anxiety and moodiness.

The Big Five Personality Traits model, or OCEAN, is fundamental in this research due to its widespread acceptance and empirical validation across various cultures and contexts. These five traits—Openness to Experience, Conscientiousness, Extraversion, Agreeableness, and Neuroticism—offer a comprehensive framework that captures the essential aspects of human behaviour and personality. This universality and depth make the Big Five particularly valuable for AI applications, as it allows for accurate and culturally sensitive personality predictions. By incorporating the Big Five into AI models, this research ensures that the analyses are

not only robust and scalable but also ethically sound, providing a solid foundation for understanding and predicting human behaviour in diverse real-world scenarios.

### 1.3.2. Myers-Briggs Type Indicator (MBTI)

The Myers-Briggs Type Indicator (MBTI) is another popular personality framework used throughout this thesis. MBTI categorizes individuals into 16 distinct personality types based on four dichotomies:

- **Extraversion (E) vs. Introversion (I):** This dichotomy reflects where individuals focus their attention and energy. Extraverts are energized by social interaction, while introverts are energized by solitary activities.
- **Sensing (S) vs. Intuition (N):** This dimension concerns how people perceive the world. Sensors focus on concrete, present information, while intuitives are more concerned with patterns and future possibilities.
- **Thinking (T) vs. Feeling (F):** This dichotomy deals with decision-making processes. Thinkers rely on logic and objective analysis, while feelers make decisions based on personal values and the impact on others.
- **Judging (J) vs. Perceiving (P):** This dimension reflects how individuals deal with the outside world. Judgers prefer structure and decisiveness, while perceivers are more flexible and open to new information.

The MBTI framework is widely used in personal development, organizational settings, and educational contexts. This research integrates MBTI to explore how deep learning models can predict and analyze these 16 personality types based on user-generated content.

### 1.3.3. Extraversion, Introversion, and Ambiversion

In addition to the Big Five and MBTI, this thesis also focuses on the specific traits of Extraversion, Introversion, and Ambiversion:

- **Extraversion:** As described in both the Big Five and MBTI models, extraversion is characterized by sociability, assertiveness, and a tendency to seek out social interactions.
- **Introversion:** In contrast, introversion is associated with a preference for solitary activities, reflection, and a lower need for social engagement.
- **Ambiversion:** Ambiversion represents a balance between extraversion and introversion. Ambiverts exhibit traits of both, depending on the context and situation. They may feel comfortable in social settings but also enjoy solitude.

The integration of personality models like the Big Five, MBTI, and the distinct traits of Extraversion, Introversion, and Ambiversion into AI research is not merely theoretical. These models serve as essential frameworks that inform the development of algorithms capable of understanding and predicting human behaviour in nuanced ways. By accurately modelling these personality traits, AI systems can provide more personalized and contextually relevant interactions, whether in educational settings, mental health applications, or personalized marketing. Moreover, the ability to detect

and analyze these traits across different languages and cultures further enhances the inclusivity and effectiveness of AI technologies, ensuring they are adaptable to a global user base. This focus on personality traits provides the foundation for the research objectives outlined in the subsequent section.

#### **1.4. Statement of Research and Research Objectives**

As we explore the complex weave of human personality, we turn to the abundant digital realm, where expressions in myriad languages offer a window into the global psyche. This exploration is driven by a fundamental research statement:

*"How can we develop and validate cross-linguistic deep learning models for personality detection that integrate psychological theories, profile emotional dispositions, and evaluate their real-world applications in domains?"*

To steer through this challenging research landscape, our study is grounded in three essential research objectives (RO):

- **RO 1: Developing and Validating Cross-Linguistic Personality Detection Models:** Our first objective is to create AI models that are not confined by language limitations but are adept at discerning personality traits from an array of linguistic contexts, offering a universal tool for personality analysis.
- **RO 2: Integrating Psychological Theories into Deep Learning Frameworks:** We aim to infuse the computational efficiency of deep learning with the qualitative depth of psychological theories. This synergy is expected to yield models that capture the nuances of personality with a degree of sophistication akin to human psychological analysis.
- **RO 3: Profiling Emotional Dispositions and Evaluating Real-World Applications:** The third objective extends beyond identification; it is about applying the insights derived from our AI models to tangible scenarios. Whether in social media, educational settings, or the workplace, we seek to demonstrate how an understanding of personality traits can significantly enhance user experience and engagement.

The core challenge addressed in this research is the development of AI models that can accurately and ethically detect and analyze personality traits from diverse and unstructured user-generated content (UGC). Traditional methods of personality assessment, such as self-reported questionnaires, are limited by biases, scalability issues, and a lack of real-time analysis. In contrast, the vast amount of UGC available today offers an opportunity to assess personality traits continuously and in real-world contexts. However, the complexity of human personality, the variability of UGC formats, and the need for cross-linguistic and cross-cultural accuracy pose significant challenges. This research seeks to overcome these challenges by integrating deep



learning techniques with established psychological theories to create models that are both precise and adaptable to diverse populations.

### **1.5. Primary Contribution of the Research**

This research advances the field of automatic personality detection through the following key contributions:

- **Development of Cross-Linguistic Models:** Introduced and validated transformer-based models tailored for English, Hindi, and Bangla, achieving high accuracy in personality detection across diverse cultural contexts.
- **Integration of Psychological Theories:** Seamlessly integrated Big Five and MBTI frameworks into deep learning models, enhancing the precision and relevance of personality trait analysis.
- **Emotional Disposition Profiling:** Developed novel methodologies for profiling emotional dispositions like optimism and pessimism, demonstrating their application in personalized AI-driven services.
- **Innovative Use of Emojis:** Pioneered the incorporation of emojis as a data source in personality prediction, providing a new dimension to text-based personality assessment.

These contributions collectively push the boundaries of AI-driven personality detection, offering robust, culturally adaptive models with practical real-world applications.

#### **1.5.1. Contribution to the Knowledge Base**

This research significantly enriches the existing knowledge base by bridging the gap between traditional psychological assessments and modern AI techniques. By introducing cross-linguistic models, it extends the applicability of personality detection across different languages and cultures, an area that has been relatively underexplored. The integration of well-established psychological frameworks with advanced deep learning architectures provides a robust, scientifically grounded approach to personality analysis, which enhances the validity and reliability of AI-driven personality predictions. Additionally, the innovative use of emojis as a predictive feature opens new avenues for understanding digital communication's role in personality expression, contributing valuable insights to both the fields of AI and psychology.

### **1.6. Organization of the Thesis**

This thesis is structured into several chapters, each focusing on a specific aspect of the research conducted. The organization of the thesis is designed to provide a



logical progression of the research process, from the introduction and literature review to the methodology, results, discussion, and conclusion.

The opening chapter, Chapter 1, introduces the burgeoning field of computational personality detection and its significance in various applications. It outlines the research objectives, which aim to develop AI models capable of understanding and classifying personality traits from user-generated content. The focus is on the integration of advanced machine learning techniques with psychological constructs to enhance personalized user experiences. This chapter underscores the innovative merger of psychological insights with deep learning illustrating the sophistication and depth these models bring to the study of personality psychology.

In chapter 2, the literature review provides a critical examination of previous works in the domain of personality detection. It explores traditional approaches and highlights the evolution toward the utilization of deep learning models. The chapter also assesses the application of psychological theories within these models, setting the stage for the research by identifying gaps and establishing the necessity for the current study. Following this, a comprehensive overview of the datasets curated for the study is presented in chapter 3. It includes a description of the multilingual datasets from various sources, the criteria for their selection, and the processes undertaken to validate their relevance and accuracy for personality analysis.

The development and validation of models for personality analysis across languages are detailed next in chapter 4. This chapter focuses on models like PersonalityBERT, which shows high accuracy in classifying personality types from English datasets, and ByaktitbaNet, a model tailored for Bangla language datasets, indicating a successful fusion of BERT embeddings with LSTM. It further introduces HindiPersonalityNet, which utilizes BioWordVec embeddings with a GRU architecture to analyze Hindi data. The section also covers the implementation of transformer models for the "Personality\_Quotes" dataset, demonstrating their versatility in detecting nuanced personality traits from text data.

Next, in chapter 5, the thesis highlights how psychological theories can be synergized with computational models to create robust frameworks for personality detection. It features the implementation of models such as KBSVE-P on Hindi datasets and discusses their performance in capturing MBTI personality types. This section also introduces the DM-MBTI Dataset Creation, leveraging the DeBERTa model to analyze decision-making styles, and discusses a predictive modeling approach applied to the Parvarish-MBTI dataset for predicting parenting styles through MBTI Profiles in the Hindi-speaking population.

Following this, the exploration into profiling emotional dispositions through transformer-based models is discussed, including their application in predicting optimism and pessimism is shown in chapter 6. The chapter also considers the practical implications of these findings in diverse settings, illustrating the real-world applicability of the research. The results and discussion chapter, chapter 7, presents a model-wise analysis of the results within each research objective. It discusses the

effectiveness and contributions of the models to the field of personality detection also addresses the limitations encountered in the research process. The final chapter, chapter 8, provides a conclusive summary of the research, outlining its significant contributions to the intersection of AI and psychology. It also delineates future research avenues, particularly the need for enhancing model interpretability and addressing ethical considerations.

Additionally, the thesis includes a list of all the scholarly articles and conference presentations that have stemmed from the research, highlighting the academic and practical contributions. A detailed compilation of all the academic and professional sources referenced throughout the thesis provides a resource for further exploration into the topics discussed.

### **1.7. Chapter Summary**

Chapter 1 introduces the complex concept of personality, emphasizing its significance across psychology, education, and AI. It discusses the limitations of traditional personality assessment methods and highlights the potential of AI, particularly deep learning and NLP models like transformers, to analyze vast amounts of user-generated content (UGC) for personality detection. This chapter also underscores the critical role of established personality models, such as the Big Five, MBTI, and traits like Extraversion, Introversion, and Ambiversion, in guiding the development of AI models that can accurately assess and predict human behaviour. Despite current research being in its infancy, there is a pressing need to develop sophisticated, cross-linguistic models that integrate psychological theories to accurately interpret human personality from digital footprints. This integration could enhance user experiences in various applications by providing more personalized and intuitive interactions, thus expanding the field's impact and utility. The chapter also outlines the research objectives, which focus on developing cross-linguistic personality detection models, integrating psychological theories with deep learning, and profiling emotional dispositions in real-world applications. Additionally, it provides an overview of the thesis organization, detailing the progression from introduction to methodology, results, discussion, and conclusion.

## CHAPTER 2

### LITERATURE REVIEW

The study of personality has long been a cornerstone of psychological research, with traditional methods such as surveys and observational techniques paving the way for understanding human behaviour. However, the advent of advanced computational technologies, particularly deep learning and data analytics, has significantly enhanced the scope and accuracy of personality detection. This section summarizes the key findings from previous research on personality detection, highlighting both traditional methods and recent advancements.

#### 2.1. Traditional Methods of Personality Detection

Historically, personality research has relied heavily on psychometric surveys and observational studies. The most widely recognized among these is the use of standardized questionnaires, such as the Myers-Briggs Type Indicator (MBTI)[16-28] and the Big Five Personality Traits[29-40] inventory. These tools are designed to measure aspects of personality that are consistent across different contexts and over time.

- **Psychometric Surveys:** These surveys ask participants to respond to a series of questions about their behaviour, preferences, and reactions to various scenarios. The responses are then analyzed to assign a personality type or score along various traits. For example, the Big Five inventory scores individuals on openness, conscientiousness, extraversion, agreeableness, and neuroticism.
- **Observational Studies:** These involve the direct observation of subjects in controlled or natural settings. Researchers look for specific behaviours that correlate with known personality traits. Such methods are often used in clinical settings or in developmental psychology to observe behavioural patterns over time.

While these traditional methods have provided valuable insights, they come with limitations. Psychometric surveys can suffer from self-reporting biases participants may answer in a way they perceive as socially desirable rather than true to their behaviour. Observational studies, on the other hand, can be time-consuming and subjective, depending on the observer's interpretations.

#### 2.2. Advancements in Personality Detection through Deep Learning and Data Analytics

The integration of machine learning, especially deep learning, into the field of personality detection represents a significant advancement, addressing many

limitations of traditional methods. Deep learning models, which use complex algorithms to process and learn from large amounts of data, can uncover patterns that are not apparent to human observers.

- **Data-Driven Insights:** Modern personality research increasingly leverages user-generated content from social media platforms as a data source. Text posts, videos, and online interactions provide a rich dataset from which to extract personality cues. For example, linguistic analysis of social media posts using text analysis tools can reveal traits such as extraversion or neuroticism based on the frequency and context of words related to emotions or social interactions.
- **Deep Learning Models:** Techniques like Natural Language Processing (NLP) have been particularly transformative. Models such as BERT (Bidirectional Encoder Representations from Transformers) [41-44] and GPT (Generative Pre-trained Transformer) [6][45] can understand the context and nuance of language used in digital communications. These models analyze text data to predict personality traits based on linguistic patterns that align with known personality dimensions.
- **Facial Recognition and Computer Vision:** Advancements in computer vision have allowed researchers to analyze video and images for micro-expressions and other non-verbal cues that indicate personality traits. Deep learning models trained on facial recognition algorithms can interpret these subtle cues, providing a non-intrusive way to assess personality.
- **Behavioural Data Analytics:** The use of behavioural data from digital footprints extends beyond text and image analysis. For instance, patterns in browsing history, app usage, and even gameplay in video games can provide insights into personality traits. Machine learning models can analyze these diverse data types to create comprehensive personality profiles.

Over the past decade, extensive research has focused on analyzing English [14-70] text data scraped from social platforms like YouTube, Facebook, and Twitter, with a significant shift towards multimodal data analysis in recent years. While the majority of studies have concentrated on English, there has been a growing interest in exploring personality detection across other languages such as Hindi [71-73], Urdu [74], Iranian [75] Indonesian [76], and Turkish [77]. This broadened linguistic scope is vital for developing more culturally inclusive AI technologies. Our research categorizes these studies into two primary tables: the first highlights work related to English datasets encompassing texts, transcripts, and audio, while the second delves into research contributions in Hindi, Urdu, Bangla [78], and Persian [79] languages.

**Table 2.1:** Literature related to Personality Detection in English Datasets

Language	Year	Author	Methodology
English	2016	Gurpinar et al. [35]	Pre-trained DCNN model was employed to extract the information using facial expressions which was combined with the scene using kernel extreme learning machine regressor. The dataset used was First Impression Chalearn.
	2016	Gurpinar et al. [36]	Used OpenSMILE tool for extracting facial emotions and acoustic features for Chalearn dataset.
	2016	Nishant Rai [46]	Over the Chalearn dataset, deep networks were used which focused on leveraging visual information.
	2016	Zhang et al. [47]	Proposed a deep bimodal regression model for Chalearn dataset.
	2017	Hernandez RK, Scott I [48]	Author worked upon MBTI Dataset using RNN model.
	2017	Ismail et al. [19]	A questionnaire analysis was made over MBTI personality from the university students based upon the courses opted.
	2017	Majumder et al. [49]	CNN-MLP model was demonstrated over Essays Dataset that is, James Pannebaker and Laura King's stream-of-consciousness essay
	2017	Gucluturk et al. [40]	Author worked upon the benchmark dataset Chalearn. They described a new web application while making prediction of the state of art model.
	2017	Miranda-Correa et al. [50]	Individual and group mood and personality recognition using short and long videos using AMIGOS dataset.
English	2018	Kampman et al. [51]	Proposed tri- modal architecture for the first impression dataset with audio as topmost layer and video as the lower most layer. Stacked CNN model was employed for each channel. The channels are fused both on decision-level and also by concatenating their respective fully connected layer.

<b>English</b>	2019	Keh and Cheng [14]	BERT model was used for personality Café forums dataset which consisted of 68K posts.
	2020	Sun et al. [52]	The author used BERTCap model over four datasets: SNIPS, StackOverFlow, ECDT and FDQuestion
	2020	Lynn et al. [53]	Data was gathered by consenting users of facebook application. Message attention and word attention with GRU model was used.
	2020	Ren et al. [54]	BERT model was used for MBTI and CNN, GRU and LSTM for essays dataset.
	2020	El-Demerdash et al. [55]	Adopted ULMFiT (Universal Language Model Fine Tuning) for text analytics task in NLP. The benchmark dataset used in this research is SoCE (Essays).
	2021	Salsabila, G. D., & Setiawan, E. B. [56]	Support Vector Machine (SVM) and BERT model was used for the data collected using BigFiveInventory questionnaire on twitter users.
	2022	Singh, S. and Singh, W. [57]	Performed a systematic literature review on the research papers from 2016 till 2022 and have categorized all the papers based upon the techniques used to detect personalities i.e. machine learning, Deep Learning, transfer learning techniques.
	2023	Ibrahim, R. T., & Ramo, F. M. [58]	Two models were proposed, first was the convolution model with five main layers while the other one was deep learning model with fuzzy learning. Both the models were tested for Big-5 personality traits.
	2023	Sirasapalli, J.J., & Malla, R.M. [59]	A conceptual framework to map MBTI personality traits with Big-5 was proposed. Later the fusion of the two datasets Essays and myPersonality was demonstrated.
	2023	Johnson, S.J., Murty, M.R. [60]	A knowledge graph-based personality detection is proposed.
	2023	Yang et al.	Proposed DeepPerson, a tool for

<b>English</b>		[61]	detecting personality from text that combines ideas from psychology and deep learning. DeepPerson uses transfer learning and attention networks along with psychological concepts and data augmentation techniques to analyze individual writing styles.
	2024	Guo et al. [62]	378 participants were asked them to fill out four personality questionnaires covering 25 personality traits, and had them perform three rounds of human machine dialogue with a pipeline task-oriented dialogue system or an end-to-end task-oriented dialogue system. We also had another 186 participants do the same with an open-domain dialogue system. The study was made using MBTI personality traits.
	2024	Grunenberg et al. [63]	This research demonstrated the feasibility of using automated predictions of personality in applied preselection and recruiting settings. CV applications of approximately 7864 candidates were uploaded to search jobs with respect to the Big-five personality traits.
	2024	Sze et al. [64]	Data of 3282 events of 144 different users engaged in activities like walking, running, etc using an iPhone sensor. Big five personality model is used to annotate the events.
	2024	Suhartono et al. [65]	250 facebook users of the myPersonality dataset was tests for five supervised learning models to determine the big-5 personality traits. Multinomial Naïve Bayes algorithm was observed to be the best one.
	2024	Hu et. al. [66]	Proposed a novel LLM for personality detection. They tested MBTI personality dimensions over the proposed model and ChatGPT to study the results.
	2024	Serrano-Guerrero et	Proposed a stacked ensemble model using multiple classifiers for personality

		al. [67]	recognition. This model was used over the gold standard myPersonality dataset.
	2024	Saeidi, S. [68]	Screenshots of the whatsapp conversations using emojis were collected and then converted to numerical value using image processing and were classified into 16 dimensions of personalities. The LSTM neural network model was proposed to determine the personality traits.
	2024	Liao et al. [69]	Authors have worked upon the two publically available datasets Chalearn and UDIVA (self-reported personality) proposed framework for audio, video and non- verbal audio-visual were designed and tested.
	2024	Alshouha et al. [70]	BERT, RoBERTa, ALBERT, ELECTRA, ERNIE, or XLNet models were tested for myPersonality dataset with 250 facebook users to detect personality.

This table 2.1 presents a comprehensive overview of the strides made in the field of personality detection using English datasets. Researchers have utilized various deep learning models and analytical techniques to extract and interpret personality traits from rich multimodal content. These studies have laid a crucial foundation for advancing our understanding of how personality is expressed and can be detected through digital interactions.

As we expand our focus beyond the English language, it becomes crucial to explore how personality detection is approached in other linguistic and cultural contexts. This exploration is essential for developing models that are not only linguistically diverse but also culturally sensitive. The following table 2.2 shifts the spotlight to significant contributions in non-English datasets, highlighting research in Hindi followed by table 2.3 showcasing the works done in Urdu, Persian, and Bangla. This shift acknowledges the rich cultural nuances that influence personality expression and detection, aiming to foster a more inclusive understanding across different global communities.



**Table 2.2:** Literature on Personality Detection in Hindi Dataset

Language	Year	Author	Methodology
<b>Hindi</b>	2013	Singh et. al. [71]	Introduced a psycho-lexical approach where the data from different Hindi novels and sources were collected and further studied based on three major Triguna such as sattvic, rajasic, and tamasic.
	2017	Singh JK and De Raad B [72]	The participation of 1250 different Hindi speaking people of young generation recorded the personality traits by observing the Triguna and the big-5 personality traits.
	2020	Khan et al. [73]	Proposed a Hindi conversational dataset "Vyaktitv" which consisted of audio, video, and Hinglish utterances (Hindi words are written using English language).

**Table 2.3:** Literature on personality detection in other non-English Languages

Language	Year	Author	Methodology
<b>Urdu</b>	2013	Khan, Iftikhar [74]	Translated English Big-five IPIP NEO to Urdu language
<b>Iranian</b>	2013	Maghsoudi Mojtaba [75]	focused on Big-five personality traits. The author wanted to understand the importance of personality for male and female bilingual learners. It was observed that bilingual female learners were extrovert than male bilingual learners.
<b>Indonesian</b>	2018	Adi et al. [76]	Presented optimization techniques for machine learning (ML)-based Bahasa Indonesian automatic personality detection.
<b>Turkish</b>	2019	Yilmaz et. al [77]	Ocean.csv dataset was translated to Turkish language. RNN was applied using LSTM. Results for test and train sets were recorded for different personality traits.
<b>Bengali</b>	2020	Rudra et. al. [78]	Used C-LSTM model for Bengali dataset available in Big-five personality traits.

<b>Persian</b>	2022	Anari et al. [79]	Proposed a lightweight deep convolutional neural network for categorizing personality into 9 types from handwritten Persian text.
----------------	------	----------------------	---

### 2.3. Major Findings from Literature Review

The comprehensive review of literature in personality analysis unveils an expansive variety of data sources and methodologies being utilized to decode human personality traits. It showcases the sophisticated integration of written texts, nonverbal behaviours, mobile data, and online gaming activities, each serving as a rich vein for extracting insightful correlations with personality traits. These findings emphasize the evolution of analytical techniques from binary classifications to complex multimodal fusion strategies, which are crucial for constructing detailed personality assessments. Despite advancements, significant gaps remain, such as biases in data collection, a narrow focus on established models like the MBTI and Big Five, and a predominance of English-centric research, all of which limit the scope and applicability of current models.

- **Diversity of Cues for Personality Analysis:**
  - i. *Written Texts:* Extensively used to infer personality traits through linguistic analysis.
  - ii. *Nonverbal Behavior:* Studies focus on gestures, facial expressions, and other nonverbal cues as important indicators of personality.
  - iii. *Mobile and Wearable Device Data:* Emerging as crucial sources for detecting behavioral patterns through user interactions and physical activity.
  - iv. *Online Games:* Player behaviors and decisions in gaming environments provide novel insights into personality traits.
  - v. *Emojis:* Increasingly recognized as valuable indicators of personality traits, given their widespread use in digital communication and ability to convey nuanced emotions and attitudes.
- **Analytical Techniques:**
  - i. *Binary Classification and Regression:* Predominantly used to correlate behavioral cues with specific personality traits.
  - ii. *Multimodal Fusion Strategies:* Being explored to integrate multiple types of data for more comprehensive personality assessments.
- **Challenges with Data:**
  - i. *Bias and Inconsistency:* Common in data collection and annotation processes, which compromise the generalizability of findings.
- **Focus on Established Models:**
  - i. *Limited Scope:* Research has primarily focused on the MBTI and Big Five personality models.

- ii. *Neglect of Other Traits:* Lesser attention to traits like introversion, extroversion, and ambiversion; which are also critical aspects of personality.
- **Linguistic Limitations:**
  - i. *English-centric Research:* Most studies utilize English language data, underscoring a need for research in other languages to accommodate global diversity.
- **Need for Comprehensive Frameworks:**
  - i. *Lack of Integrated Models:* There is an absence of a multi-faceted framework in computational psychology that encompasses thought, emotion, and behavior holistically.
  - ii. *Emotional and Attitudinal Dimensions:* Absence of models that include emotional dispositions such as optimism and pessimism, which are crucial for understanding comprehensive personality dynamics.
- **Real-World Applications:**
  - i. *Application Gaps:* Limited exploration of how personality insights apply to practical scenarios like employment and parenting.
  - ii. *Lack of Correlation Studies:* No significant research correlating personality traits with decision making, parenting styles, or real-time job application successes.
- **Demand for Interpretable Models:**
  - i. *Interpretability in Machine Learning:* There is a critical need for models that not only predict but also explain personality traits in understandable terms.

Ultimately, these findings underscore the critical need for continued innovation and expansion in the field of personality detection, to bridge the current gaps and enhance the application of this research in practical, globally diverse contexts.

## 2.4. Chapter Summary

Chapter 2 provides a comprehensive literature review on the evolution of personality detection, contrasting traditional methods with modern advancements in deep learning and data analytics. Historically, personality research has relied on psychometric surveys and observational studies, such as the Myers-Briggs Type Indicator (MBTI) and Big Five Personality Traits inventory. These methods, while valuable, suffer from limitations like self-report biases and time-consuming processes. The advent of deep learning has significantly expanded the scope of personality detection, utilizing user-generated content from social media, facial recognition, and behavioral data analytics to derive personality insights. Techniques like Natural Language Processing (NLP) with models such as BERT and GPT have

revolutionized text analysis, while computer vision advancements enable non-intrusive assessment through video and image analysis. The review highlights a shift towards multimodal data analysis, incorporating linguistic, visual, and behavioral cues to enhance personality detection. Despite progress, the research predominantly focuses on English datasets, necessitating the development of cross-linguistic models to capture cultural diversity. The chapter underscores the need for comprehensive frameworks that integrate emotional and attitudinal dimensions, address biases, and improve interpretability in machine learning models, paving the way for practical applications in employment [80-84], education [85-90], and mental health [91-95].

## CHAPTER 3

### RESEARCH DATASETS

In the realm of personality research, the integrity and diversity of data play pivotal roles in uncovering the nuanced ways individuals express and communicate their personality traits. Our study employs a carefully selected array of datasets, each offering a unique lens through which various aspects of personality are captured and analyzed. These datasets encompass a broad range of sources, from video clips and social media interactions to textual analysis and emoji [96-99] usage, reflecting the multifaceted nature of personality expression across different cultures and communication modes.

This chapter provides a comprehensive overview of the diverse datasets employed in this research, highlighting their significance in capturing the nuanced expressions of personality traits across various languages and communication mediums. The chapter emphasizes the importance of these datasets in ensuring that our analysis is both culturally inclusive and contextually relevant. By leveraging multimodal data sources—such as video clips, social media interactions, textual content, and emoji usage—this research explores the multifaceted nature of personality expression. The datasets span English, Hindi, and Bangla languages, offering a rich foundation for cross-linguistic and cross-cultural personality analysis. The overview of datasets used in this research is given in table 3.1.

**Table 3.1:** An overview of Datasets

Language	Dataset	Origin of Dataset	Dataset Labels
<i>English</i>	First Impression Chalearn	Public	Big Five Personality Traits
	Kaggle_MBTI	Public	MBTI Personality Types
	Personality_Quotes	Curated	Introvert, Extrovert, Ambivert
	DM_MBTI	Reannotated	Decision making Styles linked to MBTI Personality Types
	Emo_MBTI	Reannotated	MBTI Personality Types linked to Emoji Usage
<i>Hindi</i>	Shaksiyat	Curated	Introvert, Extrovert, Ambivert
	Vishesh_Charitr	Curated	MBTI Personality Types
	Manobhav	Reannotated	Optimistic, Pessimistic
	Parvarish	Curated	Parenting Styles linked to MBTI Personality Types
<i>Bangla</i>	Byaktitba	Curated	Introvert, Extrovert, Ambivert

### 3.1. Significance of the Diverse Spectrum of Datasets

The spectrum of datasets employed in this research is both special and important for several reasons:

- **Multimodal Data Sources:** The inclusion of video clips, social media interactions, textual analysis, and emoji usage captures the multifaceted nature of personality expression. Different modes of communication offer unique insights into how personality traits manifest in various contexts, enhancing the richness and depth of the analysis.
- **Cultural and Linguistic Diversity:** By incorporating datasets from English, Hindi, and Bangla, this study acknowledges and addresses the influence of cultural and linguistic factors on personality expression. This diversity ensures that findings are not biased towards a single cultural perspective, making the research more globally relevant and inclusive.
- **Comprehensive Personality Frameworks:** Utilizing a range of personality frameworks, including the Big Five Personality Traits, MBTI Personality Types, and other specific categorizations like introvert/extrovert/ambivert and optimistic/pessimistic outlooks, provides a holistic view of personality. This allows for a more nuanced understanding of how different personality traits interact and are perceived.
- **Reannotated and Curated Datasets:** The reannotation of certain datasets, such as DM\_MBTI and Emo\_MBTI, to highlight specific behaviours and non-verbal communication aspects, demonstrates the dynamic and evolving nature of personality research. Curated datasets ensure high-quality, contextually relevant data, enhancing the accuracy and applicability of the findings.
- **Application Across Domains:** The datasets cover various domains such as decision-making styles, parenting styles [100-102], and emotional outlooks, linking personality traits to practical aspects of daily life. This broad applicability allows the research to impact multiple fields, including psychology, education, marketing, and human-computer interaction.
- **Enhanced Analytical Robustness:** The integration of diverse datasets reduces the risk of overfitting and increases the robustness of analytical models. By validating findings across multiple data sources and types, the research ensures more reliable and generalizable results.
- **Innovation in Personality Research:** Exploring the linkage between personality traits and modern communication forms like emoji usage represents an innovative approach in personality research. This not only

keeps the study current with evolving communication trends but also opens new avenues for understanding personality in digital contexts.

- **Cross-Disciplinary Relevance:** The comprehensive and diverse nature of the datasets makes the research relevant to various disciplines beyond psychology, including linguistics, data science, cultural studies, and artificial intelligence. This cross-disciplinary relevance enhances the potential for collaborative research and practical applications.

In summary, this diverse spectrum of datasets is special and important because it provides a rich, culturally inclusive, and multi-dimensional foundation for understanding personality traits. It enhances the robustness and generalizability of the research findings, supports innovative approaches to personality analysis, and ensures relevance across multiple academic and practical domains.

In the following section, we will delve into the detailed descriptions of the datasets used in our study, focusing specifically on the English, Hindi and Bangla datasets. Understanding the origin, content, and labelling of these datasets is crucial for appreciating the depth and breadth of our research. Each dataset provides unique insights into personality traits, contributing to a comprehensive analysis of how personality is expressed and perceived across different modes of communication and cultural contexts.

### 3.2. English Language Data

The English datasets utilized in this research share several common features that make them invaluable for the study of personality traits. These commonalities ensure a robust, multifaceted approach to understanding personality expression and provide a comprehensive foundation for our analysis.

The English datasets encompass a variety of data sources, including video clips, social media content, quotes, and emoji usage. This diversity allows for a holistic examination of personality traits across different mediums and contexts, capturing both verbal and non-verbal expressions of personality. Each dataset leverages well-established personality frameworks, such as the Big Five Personality Traits and the Myers-Briggs Type Indicator (MBTI). This consistency in frameworks facilitates comparative analysis and enhances the reliability of the findings across different datasets.

The datasets include both publicly accessible data (such as those from Kaggle and Chalearn) and curated collections (like Personality\_Quotes). This combination ensures a balance between large-scale, diverse data and high-quality, contextually relevant data, enriching the overall research quality. Several datasets have been reannotated to provide deeper insights into specific aspects of personality, such as decision-making styles (DM\_MBTI) and emoji usage (Emo\_MBTI). This reannotation process adds layers of analysis, allowing for a more detailed exploration of how personality traits influence behaviour and communication.

The English datasets offer a comprehensive look at both textual and non-verbal communication. Textual datasets like Kaggle\_MBTI and Personality\_Quotes provide rich linguistic data, while non-verbal datasets like First Impression Chalearn and Emo\_MBTI capture body language and emoji usage, respectively. This dual focus ensures a thorough understanding of personality expression. By including modern forms of communication, such as social media interactions and emoji usage, these datasets remain highly relevant to contemporary communication trends. This relevance ensures that the research findings are applicable to current and evolving modes of human interaction. Although centred on English language data, these datasets likely encompass a wide range of cultural and demographic backgrounds, given the global reach of platforms like social media and public competitions. This variety supports a more inclusive and representative analysis of personality traits.

Undoubtedly, the English datasets in this research are characterized by their diverse data sources, use of established personality frameworks, accessibility, enhanced insights through reannotation, comprehensive analysis of communication modes, relevance to modern trends, and cultural diversity. These common features collectively contribute to a robust and nuanced understanding of personality traits, making the datasets integral to our study.

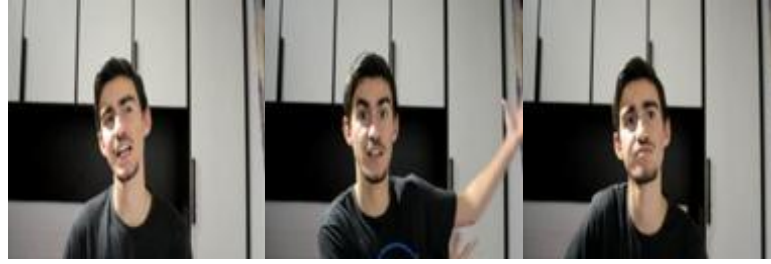
### 3.2.1. First Impression Chalearn (English, Public)

This dataset is comprised of video clips where annotators assess subjects based on the Big Five Personality Traits: openness, conscientiousness, extraversion, agreeableness, and neuroticism. It provides a visual and dynamic medium to study personality traits, allowing for the examination of non-verbal cues and body language. It is constituted of 10,000 videos that have an average duration of 15 seconds, extracted from more than 3,000 HD Youtube videos of people speaking in English while facing the camera. Each video has five labels in the range [0,1] which depict the big five personality traits. The videos are also accompanied by their transcripts, in the form of a pickle file. The transcript contains 435984 words, averaging to 43 words per video, with a maximum of 85 words in a single transcript and 14535 unique words.

0	He's cutting it and then turn around and see t...	0.523364	0.488889	0.626374	0.552083	0.601942
1	Responsibility to house the organ I had been g...	0.345794	0.366667	0.472527	0.375000	0.582524
2	I actually got quite a few sets of black pens ...	0.252336	0.511111	0.406593	0.291667	0.485437
3	I ate a lot. I'd like a lot of foods. I rememb...	0.457944	0.377778	0.505495	0.489583	0.398058
4	Now I'll ask you guys to leave a question in t...	0.607477	0.622222	0.406593	0.489583	0.621359
...	...	...	...	...	...	...
9995	Du du du. No, that's not why I made this video...	0.289720	0.300000	0.208791	0.312500	0.135922
9996	They do it all the time.	0.719626	0.722222	0.670330	0.781250	0.572816
9997	Comfortable and I don't want anyone to feel un...	0.355140	0.677778	0.472527	0.395833	0.446602
9998	You're not giving yourself enough calories to ...	0.467290	0.622222	0.527473	0.645833	0.669903
9999	Eat enough carbs, eat more fats to get in more...	0.654206	0.588889	0.813187	0.635417	0.728155

**Fig. 3.1.** Transcripts of the dataset



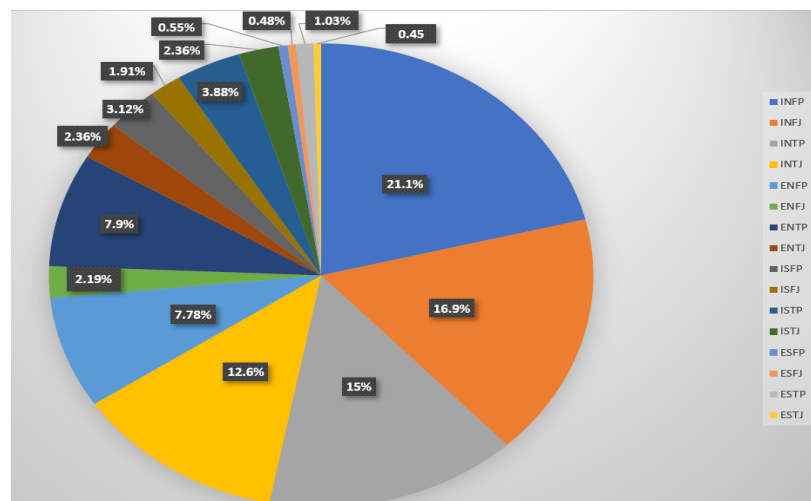


**Fig.3.2.** Frames of the Dataset

### 3.2.2. Kaggle\_MBTI (English, Public)

Sourced from Kaggle, this dataset categorizes individuals according to the Myers-Briggs Type Indicator (MBTI), offering insights into various personality types. This dataset is valuable for analyzing written expressions and social media content to understand personality distribution and interactions among different MBTI types.

The concept of Kaggle\_MBTI is that every persona consists of four magnitudes, and each magnitude has 2 possibilities [31]. These dimensions or the magnitudes and their corresponding possibilities are Introversion (I) Extroversion (E), Intuition (N) Sensing (S), Thinking (T) Feeling (F), and Judging (J) Perceiving (P). In total, this dataset has 8675 rows along with 16 combinations or personality types, which are denoted using a four-letter abbreviation. It means that each person will have a personality abbreviated such as ENFJ, INFJ, ENFP, INFP, ENTJ, INTJ, ENTP, INTP, ESFJ, ISFJ, ESFP, ISFP, ESTJ, ISTJ, ESTP, and ISTP from the permutation of all the four axis. For example, someone who is an extrovert, relies more on sensing, thinking, and perceiving rather than judging and will be labelled as an ESTP. As far as the data in the Kaggle\_MBTI is concerned, in a single row, it consists of around 50 social media posts of a specific user separated by "|||". The dataset has no null or missing values neither in the type of attribute nor in the post attribute. The following fig. 3.3 shows the personality traits distribution in the Kaggle\_MBTI dataset.



**Fig. 3.3.** Personalities distribution in the MBTI dataset

### 3.2.3. Personality\_Quotes (English, Curated)

This curated dataset includes quotes and statements categorized into introverts, extroverts, and ambiverts. It aids in textual analysis of personality expression, allowing researchers to explore how different personality types articulate their thoughts and emotions through language.

The development of the *Personality\_Quotes* dataset involved the application of web scraping techniques to aggregate user-generated quotes associated with three distinct personality archetypes: introvert, extrovert, and ambivert. Each personality category comprises approximately 350 textual quotes. These quotes were systematically generated de novo through the utilization of web scraping methodologies, involving the extraction of textual content from diverse online platforms. A comprehensive tabular representation of the dataset's statistical particulars is provided in Table 3.2.

**Table 3.2.** *Personality\_Quotes* Dataset statistics

Description	Statistics
Total number of instances	1028
Total number of Extrovert instances	365
Total number of Introvert instances	385
Total number of Ambivert instances	278

The dataset constitutes 1029 instances organized into a structured format with two columns. The first row of the dataset is used for the header, which describes the attributes present in the dataset. The main column contains textual quotes, while the second column contains categorical labels that classify each instance as either an extrovert, ambivert, or introvert expression. A visual representation, shown in fig.3.4, displays these quotes alongside their corresponding categorical labels. Regarding the categorical labels, label 0 is used to represent instances displaying introverted traits, label 1 corresponds to instances showcasing extroverted attributes, and label 2 is assigned to instances exemplifying ambivert characteristics.

	text	label
576	extroverts want us to have fun because they as...	1
542	they say that extroverts are unhappier than in...	1
844	my life is a constant battle between needing ...	2
102	introvert fact we love our alone time	0
257	stay true to your own nature if you like to do...	0
...	...	...
534	as an introvert you can be your own best frien...	1
52	inside was where she lived physically and ment...	0
998	when it comes to trusting other people someti...	2
996	some people think i am quiet while others thi...	2
932	my personality confuses people i enjoy being a...	2

**Fig.3.4.** *Personality\_Quotes* Dataset snapshot

The *Personality\_Quotes* dataset introduced in this context holds great potential as a foundational asset for creating a benchmark in the field of personality detection within user-generated content. Notably, the abundance of quotes found in posts shared across diverse social media platforms like Facebook and Instagram makes this dataset particularly well-suited for driving advancements in research within this field. The use of introversion, extroversion, and ambiversion as personality dimensions in this study is justified for several reasons. Firstly, these traits are well-established and widely recognized personality dimensions in psychological research. The introversion-extroversion spectrum captures fundamental variations in how individuals respond to social stimuli, with introverts generally being more reserved and reflective, and extroverts being more outgoing and sociable. The addition of ambiversion acknowledges the nuanced nature of personality, encompassing individuals who exhibit characteristics of both introversion and extroversion.

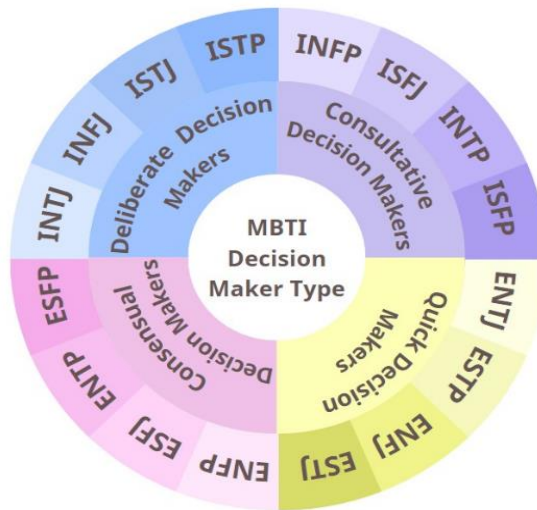
Furthermore, introversion, extroversion, and ambiversion traits are accessible and easily understood by both researchers and the general public. This simplifies the annotation and interpretation of personality labels in the dataset, promoting consistency and facilitating communication about the research findings. Additionally, these dimensions have real-world implications, influencing various aspects of an individual's life, from social interactions to career choices. By focusing on introversion, extroversion, and ambiversion, the study addresses a relatable and relevant aspect of human behaviour, contributing to the potential practical applications of the research. Overall, the choice to use introversion, extroversion, and ambiversion as personality dimensions is grounded in their established psychological significance, ease of interpretation, and their relevance in understanding and characterizing individual differences in personality expression.

#### **3.2.4. DM\_MBTI (English, Reannotated)**

Initially focused on MBTI types, this dataset has been reannotated to highlight decision-making styles linked to these personality types. This offers a unique perspective on behavior analysis, providing insights into how different personality types of approach decision-making processes.

The Kaggle\_MBTI dataset is structured around the concept that each individual's personality is composed of four dimensions, each with two potential attributes. These dimensions and their corresponding attributes include Introversion (I)Extroversion (E), Intuition (N)Sensing (S), Thinking (T)Feeling (F), and Judging (J)Perceiving (P). In total, this dataset comprises 8,675 rows, representing all 16 possible combinations of these personality dimensions. These combinations are denoted by unique four-letter abbreviations, such as ENFJ, INFJ, ENFP, INFP, ENTJ, INTJ, ENTP, INTP, ESFJ, ISFJ, ESNP, ISNP, ESTJ, ISTJ, ESTP, and ISTP, which emerge from the various permutations of the four dimensions. In this study, we have relabelled the dataset with a new category called "Decision Maker Type". To achieve this, we have identified four broad decision-making categories, as illustrated in Fig. 3.5:

- **Deliberate Decision Makers:** These individuals exhibit decision-making styles that align with specific MBTI types, including INTJs, INFJs, ISTJs, and ISTPs. Deliberate decision makers are likely to approach choices thoughtfully and systematically. They often take time to consider all aspects of a situation, preferring to analyze and weigh options systematically before arriving at a conclusion. This approach indicates a preference for thoroughness and precision in decision-making.
- **Consultative Decision Makers:** Consisting of MBTI types INFP, ISFJ, INTP, and ISFP, this group tends to value input from others and considers various perspectives before making decisions. These individuals are more inclined to engage in a collaborative decision-making process, where they seek advice and opinions from different sources to ensure a well-rounded understanding of the issue at hand. Their decision-making style is marked by a desire to understand and incorporate diverse viewpoints.
- **Quick Decision Makers:** These individuals are matched with MBTI types such as ENTJs, ESTPs, ENFJs, and ESTJs. Individuals in this group are known for their swift and confident decision-making. They often rely on their intuition and a sense of decisiveness, which allows them to make quick judgments. This style is beneficial in fast-paced or high-pressure environments where rapid decisions are necessary. These types are often seen as bold and assertive in their decision-making approach.
- **Consensual Decision Makers:** This category is associated with MBTI types like ENFPs, ESFJs, ENTPs, and ESFPs. Characterized by a focus on collaboration and consensus, these individuals prioritize group harmony and collective agreement in their decision-making processes. They tend to involve others extensively and strive to reach decisions that are acceptable and satisfactory to all involved parties. Their approach is often democratic, valuing the input and agreement of the group over individual preferences.



**Fig.3.5.** Decision maker-MBTI Personality Trait Mapping

By relabelling the MBTI types into these four decision-making categories, the DM-MBTI dataset offers a fresh perspective on how personality types can influence an individual's approach to making decisions. It underscores the diversity in decision-making styles and the importance of understanding these differences in various contexts, such as team dynamics, leadership, and personal growth. Fig.3.6 visually presents the posts from the dataset, each labelled with one of the four decision-maker types: Consensual, Consultative, Deliberate, and Quick.

DecisionMakerType	PersonalityType	Posts
Consensual	ENTP	im finding the lack of me in these posts very alarming sex can be boring if its in the same position often for example me and
Consultative	INTP	good one of course to which i say i know thats my blessing and my curse does being absolutely positive that you and your
Quick	ENTJ	youre fired thats another silly misconception that approaching is logically is going to be the key to unlocking whatever it is you
Deliberate	INTJ	i tend to build up a collection of things on my desktop that i use frequently and then move them into a folder called everythin
Consultative	INTP	im in this position where i have to actually let go of the person due to a various reasons unfortunately im having trouble must
Deliberate	INFJ	one time my parents were fighting over my dads affair and my dad pushed my mom the fall broke her finger shes pointed a
Deliberate	INTJ	fair enough if thats how you want to look at it like i stated before they were incredibly naive in their comments however they
Consultative	INTP	basically this can i has cheezburgr i am very fond of my top hat too i certainly did not expect to see a thread about top hats

**Fig.3.6.** Snapshot of DM-MBTI Dataset

### 3.2.5. Emo\_MBTI (English, Reannotated)

This reannotated dataset links MBTI personality types with emoji usage, exploring non-verbal communication aspects of personality. By examining how individuals use emojis to express their emotions and personality traits, this dataset provides a modern take on personality analysis in digital communication.

The primary resource for this was the MBTI Dataset, which comprises snippets from the last 50 posts of over 8600 individuals, each segmented and categorized into one of the 16 Myers-Briggs Type Indicator (MBTI) personality types. Entries within this dataset are delineated by "|||" (three pipe characters), organizing the text to reflect distinct posting instances by each user. To enrich the foundational MBTI Dataset, we employed advanced web scraping techniques

targeting the r/mbti subreddit on Reddit. Our focus was on users who frequently use emojis in their posts, a criterion that suggests heightened emotional expression. Specifically, we collected the last 50 posts from users who utilized more than three emojis per post, using these emojis as contextual labels for the text. Using the Python Reddit API Wrapper (PRAW), we meticulously extracted data from users on the r/mbti subreddit. The initial scraping yielded 6280 data points, which were then filtered focus on users demonstrating consistent engagement:

- 6029 data points included users with more than 10 comments.
- 5527 data points were further refined to include only those with more than 50 comments, ensuring a robust representation of active community members.

	label	text
index		
0	🤔 🤔 🤔 🤔	While currently on trial for fraud     The Sau...
1	🤔 🤔 🤔 🤔	I love how all these analysts and fans contin...
2	🤔 🤔 🤔 🤔	I immediately burst into tears     Look into ...
3	🤔 🤔 🤔 🤔	No     sorry     I had b rob too it feels b...
4	🤔 🤔 🤔 🤔	Thanks man     I may hit you up I'm in SEPA...
5	🤔 🤔 🤔 🤔	Rock vs Austin with Mankind as special ref    ...
6	🤔 🤔 🤔 🤔	That makes total sense     Lolz     the ide...
7	🤔 🤔 🤔 🤔	Heidi Gardner is "FORTY"     Revoke Stefan Kra...
8	🤔 🤔 🤔 🤔	That still seems really low but I don't know w...
9	🤔 🤔 🤔 🤔	Some books I have used and like are* Suomen m...

**Fig.3.7.** MBTI Posts relabelled with Emojis

This process not only expanded our dataset but also ensured the relevance and depth of the data collected, focusing on users with significant subreddit participation. Post-collection, the data underwent processing where a sophisticated model was trained to map text to emojis effectively. This model leverages the nuances of sequence-to-sequence generation, using generative AI models, specifically Large Language Models (LLMs), known for their capability in handling complex language tasks. For the task of mapping text to emojis, three sophisticated models were utilized: FlanT5, Pegasus, and BART, each chosen for their unique strengths in summarization and contextual understanding.

- **FlanT5:** This transformer-based model excels at generating concise, contextually accurate summaries by deeply understanding the intricacies within extensive texts. Its proficiency in handling complex narratives makes it highly effective for detailed summarization tasks.
- **Pegasus:** Renowned for its abstractive summarization capabilities, Pegasus employs a unique pre-training strategy tailored specifically for summarization. This approach enables it to distil extensive text into essence-focused summaries, adeptly capturing the core messages.



- **BART:** Utilizing a denoising autoencoder architecture, BART is designed to reconstruct and simplify text, making it particularly suitable for summarizing complex content into more digestible forms, including the translation of text sentiments and themes into corresponding emojis.

These models were chosen for their ability to efficiently process and condense textual information, facilitating accurate and expressive emoji mappings that reflect the nuanced emotional and thematic undertones of the original content. The dataset after mapping is as follows:

Unnamed: 0	type	posts	emoji
0	INFJ	'http://www.youtube.com/watch?v=qsXHcwe3krw   ...	👤 📺 📺 📺 📺 📺
1	ENTP	'I'm finding the lack of me in these posts ver...	👤 🤔 🤔 🤔
2	INTP	'Good one _____ https://www.youtube.com/wat...	💡 📺 📺 📺 📺 📺
3	INTJ	'Dear INTP, I enjoyed our conversation the o...	👤 📺 📺 📺 📺 📺
4	ENTJ	'You're fired.    That's another silly misconce...	👤 📺 📺 📺 📺 📺
...	...	...	...
8670	ISFP	'https://www.youtube.com/watch?v=t8edHB_h908   ...	👤 📺 📺 📺 📺 📺
8671	ENFP	'So...if this thread already exists someplace ...	👤 📺 📺 📺 📺 📺
8672	INTP	'So many questions when i do these things. I ...	👤 📺 📺 📺 📺 📺
8673	INFP	'I am very conflicted right now when it comes ...	👤 📺 📺 📺 📺 📺
8674	INFP	'It has been too long since I have been on per...	👤 📺 📺 📺 📺 📺

**Fig.3.8.** Snapshot of Emoji Labelled MBTI Dataset

This enriched dataset, supplemented with emoji labels, offers novel insights into personality types through the lens of textual and emotive expression.

### 3.3. Hindi Language Data

In this section, we explore a series of datasets curated from Hindi language sources, each designed to analyze different aspects of personality traits within Hindi-speaking populations. These datasets enable researchers to understand how personality traits are expressed and perceived among speakers of Hindi, providing a rich foundation for cross-cultural studies in personality psychology.

The Hindi language, with its complex structure and diverse dialects, presents unique challenges and opportunities for natural language processing (NLP). By curating and annotating dialogues and texts from various sources, we aim to create robust models that can accurately identify and classify personality traits based on linguistic cues. Each dataset discussed in this section has been meticulously prepared to ensure high-quality annotations and to represent a wide range of expressions and sentiments. From transcribed dialogues of popular Hindi series to annotations based on the MBTI personality types, these datasets are tailored to support a variety of NLP tasks, including sentiment analysis, personality detection, and more.

### 3.3.1. Shaksiyat (Hindi, Curated)

A curated dataset in Hindi, Shaksiyat classifies individuals into introverts, extroverts, and ambiverts. This dataset is instrumental in studying personality traits within the context of Hindi-speaking populations, allowing for cross-cultural comparisons.

To create the शख्सियत (pronounced as Shakhshiyat) dataset, which means "personality" in Hindi, we transcribed the dialogues of the Indian crime-thriller drama series called Aarya in Hindi. We manually classified the dialogues into three personality categories: introvert, extrovert, and ambivert, using various differentiators and context to guide our classifications. In order to ensure the accuracy of the annotations, we conducted a validation process in collaboration with a team of two highly knowledgeable and qualified post-graduate students specializing in Psychology from a reputable public university located in Delhi, India. Table 3.3 provides statistical details of the dataset.

**Table 3.3.** Shakhshiyat (शख्सियत) Dataset statistics

Description	Statistics
Total number of instances	6734
Total number of “Extrovert” instances	1710
Total number of “Introvert” instances	495
Total number of Ambivert instances	4529

The dataset comprises 6734 rows and 3 columns, with the first row serving as the header. The first column contains the speaker's name, such as Aarya, Veer, Tej, Adi, etc. The second column contains the actual utterance spoken by the speaker, while the last column provides the label for the utterance (extrovert, ambivert, or introvert). Table 3.4 illustrates some sample dialogues along with their corresponding categories. Our proposed Shakhshiyat (शख्सियत) dataset can serve as a benchmark for personality detection in the Hindi language, as well as other text classification or NLP-based tasks related to Hindi language.

**Table 3.4.** Snapshot of Shakhshiyat (शख्सियत) dataset

Speaker (Character)	Utterances (Hindi)	Label
AARYA	हाय स्वीटी।	Extrovert
	Hi sweetie.	
ADI	-माँ, वीर मुझे परेशान कर रहा है!	Introvert
	Mother, Veer is troubling me!	
AARYA	मैं जानता हूँ। घबराओ मत।	Ambivert
	I know. Do not panic.	
AARYA	घबराना बंद करो।	Introvert
	stop panicking	
VEER	आदि, बैठो।	Ambivert



	Adi, sit down.	
AARYA	आप अभी उठे हो	Extrovert
	you just woke up	
VEER	पापा, उससे कहो कि मुझे दे दो	Introvert
	papa, tell him to give it to me	
VEER	वहाँ बैठो!	Ambivert
	sit there!	
ARU	माँ, मुझे पैसे चाहिए।	Extrovert
	Mother, I need money.	

The data extracted from the transcripts of the Hindi series contains various special characters, URLs, etc. This part of the transcript adds non-pertinent noise to the dialogues, and they need to be removed. To tackle this, we perform various text-cleaning techniques used extensively in Natural Language Processing. Since our dataset is in the Hindi language, we also use iNLTK library to better process the obtained Hindi dialogues. The iNLTK library is publicly available with basic in-built functions for Natural Language Processing in Indian Languages.

### 3.3.2. Vishesh\_Charitr (Hindi, Curated)

The "Vishesh Charitr" dataset focuses on the MBTI personality types, emphasizing the study of personality traits within Hindi-speaking communities. This dataset facilitates the understanding of how different MBTI types are perceived and articulated among Hindi speakers.

This dataset, known as "विशेष चरित्र\_MBTI," is a specialized low-resource Hindi language collection. It compiles and annotates the utterances (quotes) of various individuals found on Google using the Google API. The dataset comprises 1824 annotated entries, each labeled with one of the 16 MBTI personality type abbreviations. These labels were determined with the assistance of a psychology expert. Fig.3.9 displays the annotated examples from the first five rows of the "विशेष चरित्र\_MBTI" dataset.

	type	posts	IE	NS	TF	JP
0	ENFJ	1. "अगर एक आदमी सब कुछ नष्ट कर सकता है, तो एक ...	0	1	0	1
1	ENFJ	2. "नुकसान की गारंटी देने का सबसे अच्छा तरीका ...	0	1	0	1
2	ENFJ	3. "आप जिस सड़क पर हैं उसे बदलने का अभी भी समय...	0	1	0	1
3	ENFJ	4. "डर से झिझक होती है, और झिझक आपके सबसे बुरे...	0	1	0	1
4	ENFJ	5. "मुझे कोई पछतावा नहीं है। यदि आप चीजों पर प...	0	1	0	1

**Fig.3.9.** Sample of Hindi dataset विशेष चरित्र\_MBTI

### 3.3.3. Manobhav (Hindi, Reannotated)

The development of the Hindi dataset ‘मनोभाव’ (pronounced as Manobhav) included re-annotating an existing Indian low-resource Hindi language dataset ‘विशेष चरित्र\_MBTI’ (pronounced as vishesh charitr) dataset with pessimistic and optimistic class labels for profiling text-based emotional attitude. The primary dataset has a total of 969 quotes scrapped through Google and labelled for 16 different personality traits of MBTI trait theory as INTP, ISTP, INTJ, ISTJ, INFP, ISFP, INFJ, ISFJ, ENTP, ESTP, ENTJ, ESTJ, ENFP, ESFP, ENFJ and ESFJ. It is now with the help of psychologists were relabelled and mapped these traits of human personality into optimistic and pessimistic emotional attitudes that individuals hold toward various situations, events, and outcome. The ‘मनोभाव’ (pronounced as Manobhav) dataset has 969 quotes out of which 488 quotes are labelled as optimistic while the rest 481 are labelled as Pessimistic. The snapshot is shown in fig.3.10.

	type_index	text
0	Optimistic	अगर एक आदमी सब कुछ नष्ट कर सकता है, तो एक लड़क...
1	Optimistic	नुकसान की गारंटी देने का सबसे अच्छा तरीका छोड़...
2	Optimistic	आप जिस सड़क पर हैं उसे बदलने का अभी भी समय है
3	Optimistic	डर से झिझक होती है, और झिझक आपके सबसे बुरे डर ...
4	Optimistic	मुझे कोई पछतावा नहीं है। यदि आप चीजों पर पछताव...
...	...	...
1438	Optimistic	हेलो *ईएनटीपी ग्रिन* बस इतना ही चाहिए। जब तक म...

**Fig.3.10** ‘मनोभाव’ Dataset snapshot

While the MBTI doesn't explicitly measure traits like pessimism and optimism, certain aspects of personality types might be associated with these emotional attitudes. A general overview of how certain MBTI personality types might relate with optimism and pessimism are as follows:

- **Optimistic Attitude:**

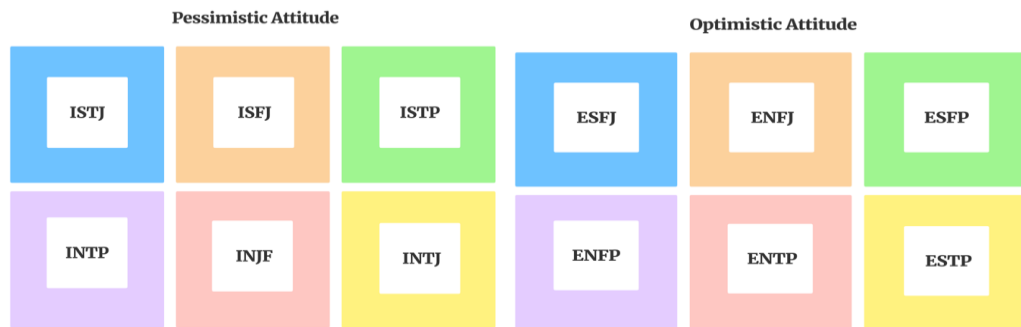
- i. Extraverted Intuitive Types (ENFP, ENTP): These types like brainstorming and are open to new ideas. They might have an optimistic outlook due to their focus on potential and their tendency to see multiple angles.
- ii. Extraverted Feeling Types (ESFJ, ENFJ): Valuing harmony and relationships, these types may lean toward optimism as they prioritize positive interactions and often work to uplift others.
- iii. Perceiving Types (ESFP, ESTP, ENFP, ENTP): Perceiving types are more adaptable and open-ended, which can lead to an optimistic outlook

since they are comfortable exploring different options and adjusting to changing circumstances.

- **Pessimistic Attitude:**

- Introverted Sensing Types (ISTJ, ISFJ): These types often focus on past experiences and known information. They might exhibit cautiousness, which can sometimes lead to a more pessimistic viewpoint when dealing with uncertainty.
- Introverted Thinking Types (ISTP, INTP): Their analytical and critical thinking can sometimes lead them to consider potential problems and challenges, which might align with a more cautious or pessimistic mindset.
- Judging Types (ISTJ, ISFJ, INFJ, INTJ): Judging types value structure and organization, which can sometimes lead to a more realistic or cautious approach, potentially leaning toward pessimism.

The following fig.3.11 visually summarizes this mapping.



**Fig. 3.11.** Mapping MBTI personality Traits to Pessimistic-Optimistic Attitudes

#### 3.3.4. Parvarish (Hindi, Curated)

A unique dataset that links parenting styles to MBTI personality types, Parvarish offers a comprehensive look at personality development influenced by parenting. It helps in understanding the interplay between parental influence and personality traits in Hindi-speaking populations.

This hindi dataset Parvarish (परवरिश) is based on two parenting styles: Energetic Explorers (EE) and Steadfast Guardians (SG). The data is binary labelled where '0' is for energetic explorers while '1' for steadfast guardians. Fig. below shows the snapshot of the dataset. The dataset consists of 1000 posts of hindi language. The snapshot is shown in the following fig. 3.12.

MBTI_Type	Posts	Parenting_Style
ENFJ	"अगर एक आदमी सब कुछ नष्ट कर सकता है, तो एक लड़की इसे क्यों नहीं बदल सकती?"	Energetic Explorer
ENFP	"यदि आप महानता प्राप्त करना चाहते हैं, तो अनुमति मांगना बंद करें।"	Energetic Explorer
ENFP	"जीवन में रुचि विकसित करें जैसा कि आप इसे देखते हैं; लोगों में, चीजों में, साहित्य में, संगीत में - दुनिया इतनी समृद्ध है, बस समृद्ध खजाने, सुंदर आत्माओं और दिलचस्प लोगों के साथ धड़कता है। अपने आप को भूल जाओ। "	Energetic Explorer
ENTJ	हम इस बारे में क्या बात कर रहे हैं?	Steadfast Guardian
ENTP	"प्रश्न हमेशा उत्तर से अधिक महत्वपूर्ण होते हैं।"	Energetic Explorer
INFJ	"हर दिन एक ऐसा काम करें जिससे आपको डर लगे।"	Steadfast Guardian
INFP	"मैं एक गहरा सतही व्यक्ति हूँ।"	Steadfast Guardian
INTJ	"मेरी शक्तियाँ सामान्य हैं। केवल मेरे आवेदन से ही मुझे सफलता मिलती है।"	Steadfast Guardian

**Fig.3.12.** Snapshot of Parvarish dataset

### 3.4. Bangla Language Data: Byaktitba (ব্যক্তিত্ব)

For creating the Byaktitba (ব্যক্তিত্ব, meaning personality in Bengali) dataset, we have transcribed the dialogues of the Bengali version of the Indian epic action film, Baahubali 2: The Conclusion (बाहुबलि २: द्या कनक्किउशान) and manually annotated the dialogues into three personality categories, namely, introvert, extrovert and ambivert. Various differentiators and context help in classifying the dialogues or utterances into their respective categories. The annotations were further validated with the help of a team of 3 Psychology post-graduates from a public-funded university in Delhi, India. The statistical detail of the dataset is presented in table 3.4.

**Table 3.4.** Byaktitba (ব্যক্তিত্ব) Dataset statistics

Description	Statistics
Total number of instances	1476
Total number of "Extrovert" instances	479
Total number of "Introvert" instances	37
Total number of "Ambivert" instances	960

The dataset consists of 1477 rows and 3 columns. The first row is the heading itself, whereas the first column depicts the speaker (the characters: Baahubali, Shivsena, Devsena etc.), the second column depicts the utterance itself, and the last column depicts the label for the utterance (extrovert, ambivert, introvert). Fig. 3.13 shows the sample dialogues along with their respective categories. The Byaktitba (ব্যক্তিত্ব) Dataset proposed by us can be used as a benchmark dataset for personality detection in the Bengali language or for other text classification or NLP-based tasks related to the Bengali language.

Speaker (Character)	Utterance Bengali	Personality
Kattappa	যুদ্ধের সঙ্গে কালকেয়াস ধরে ছিল, Calcaeus was holding on with the war	Ambivert
Avantika	স্বপ্নরাজ্য আপনাকে আমন্ত্রণ আনন্দিতভাবে The kingdom of dreams happily invites you	Extrovert
Kattappa	আমি পালিয়ে মৃত্যুর দ্বারা একটি চুল এর প্রসার! I escaped death by a hair's breadth!	Introvert
Sivagami	সহচর, সবার জন্য একটি নিরাপদ স্থান Companion, a safe place for everyone	Ambivert
Baahubali	অস্ত্র weapons	Extrovert
Kattappa	আর্মি একা করার জন্য যথেষ্ট না হয় তাদের বিরোধিতা Army alone is not enough to oppose them	Introvert
Baahubali	আজ যখন সে উইল করা সিংহাসনে বোঝানো করার জন্য আপনি যে নিম্ন জীবন Today when she is the lower life that you mean to the throne to be will	Ambivert
Avantika	আপনি কোথায় যাচ্ছেন? Where are you going?	Extrovert
Deva Sena	কান্না, ঘুম নির্দোষভাবে Cry, Sleep Innocently	Introvert

**Fig.3.13** Snapshot of Byaktitba (ব্যক্তিত্ব) Dataset

### 3.5. Validation and Testing of Datasets

Our research methodology is thoughtfully designed to ensure a comprehensive analysis of personality traits across a variety of datasets, utilizing a blend of advanced deep learning and transformer models, including the cutting-edge large language models (LLMs). The upcoming section provides a detailed overview of the datasets employed in our study, followed by an exploration of how each research objective leverages these sophisticated technologies to advance the field of personality analysis. This structured approach allows us to integrate complex computational techniques with deep psychological insights, aiming to revolutionize our understanding and application of personality assessments.

To ensure the highest standards of data accuracy and relevance, each dataset was meticulously chosen and subjected to a rigorous extraction and validation process. This comprehensive approach was essential to maintain the integrity and applicability of the data to the study's objectives. The validation process was multifaceted, involving several critical steps to guarantee the robustness and scientific rigor of the datasets.

Initially, the selection of datasets was guided by their direct relevance to the research objectives, focusing on those that best captured the nuances of personality traits across various mediums and cultural contexts. This careful selection was the

first step in ensuring that the data would be pertinent and useful for the intended analyses. Once selected, the datasets underwent a meticulous extraction process. This phase involved a detailed review and curation of the data to eliminate any noise or irrelevant information that could compromise the validity of the research. Professional psychologists played a crucial role in this step, leveraging their expertise to ascertain the psychological validity of the data. Their involvement was pivotal in ensuring that the datasets accurately reflected the psychological constructs they were intended to measure.

To further enhance the reliability of the datasets, a sophisticated validation framework was implemented, incorporating advanced interactions with OpenAI's language model, ChatGPT. This model was utilized to verify the consistency and reliability of the psychological constructs within the datasets. By engaging with ChatGPT, the research team could simulate and test various scenarios and responses, ensuring that the data's psychological representations were coherent and reliable.

This rigorous validation process was not a one-time procedure but an ongoing effort to continuously refine and improve the datasets. Feedback loops were established, where insights from preliminary analyses were used to further scrutinize and enhance the data quality. This iterative approach ensured that any potential discrepancies were promptly addressed, maintaining the highest standards of data integrity throughout the research.

The combination of expert psychological oversight and advanced AI validation ensured that the datasets were not only diverse and robust but also aligned with the scientific rigor required for advanced psychological research. This thorough validation and testing framework provides a solid foundation for the study, ensuring that the findings derived from these datasets are both reliable and valid.

In summary, the validation and testing of the datasets involved a meticulous selection process, expert psychological review, advanced AI interaction for reliability checks, and continuous refinement. This comprehensive approach guarantees that our datasets meet the stringent requirements necessary for high-quality psychological research, thereby enhancing the credibility and impact of the study.

### **3.6. Preprocessing of Datasets**

The preprocessing of datasets is a critical step in preparing the data for analysis, ensuring that it is clean, consistent, and suitable for use in deep learning models. This section outlines the specific preprocessing steps applied to the various datasets used in this research, with particular attention to the Hindi and Bangla datasets, as well as the handling of emojis, audio, and video data.

#### **3.6.1 Text Cleaning and Normalization**

For all datasets, including those in Hindi and Bangla, text cleaning was performed to remove noise and irrelevant information. The key steps included:

- **Removal of Symbols:** We methodically eliminated punctuation marks, hyperlinks, and specialized symbols from the text in all datasets. This step is



essential for reducing noise in the data, which can otherwise lead to inaccuracies in analyses. Characters such as "@", "#", and URLs can distract models from learning meaningful patterns, so their removal is crucial for enhancing data clarity.

- **Tokenization:** The text was fragmented into separate linguistic units or tokens. Tokenization allows for the analysis of individual words or phrases, which is a foundational step in most natural language processing (NLP) tasks. It helps in converting a bulk of text into smaller, manageable pieces for subsequent processing.
- **Elimination of Stop Words:** We removed commonly used but insignificant words (stop words), such as "the", "is", and "and". This step enhances the linguistic integrity of the dataset by focusing on more meaningful words that contribute to the semantic understanding of the text. Removing stop words reduces the dimensionality of the data and improves processing efficiency.
- **Standardization through Stemming and Lemmatization:** We standardized word forms using stemming and lemmatization procedures. This process is crucial for achieving linguistic congruence across different forms of the same word, aiding in the consistency of textual analysis. It helps to group similar words, enhancing the model's ability to find relevant patterns.
- **Cleaning Special Characters and URLs:** The datasets contained various special characters and URLs, which added non-pertinent noise to the data. Removing these elements is vital as they can skew the analysis and lead to misleading results. We utilized text-cleaning techniques and NLP libraries to process the dialogues and text effectively, ensuring cleaner and more focused input data.

### 3.6.2 Handling of Emojis

Emojis were recognized as significant elements within the textual data, particularly in the context of personality detection. The steps involved in processing emojis were as follows:

- **Emoji Tokenization:** Emojis were treated as unique tokens and integrated into the text sequences. This allowed the models to learn the contextual significance of emojis alongside regular text.
- **Contextual Analysis:** The meaning of emojis was considered within the context of the surrounding text and cultural nuances. This analysis ensured that the models could interpret the emotional tone conveyed by emojis, which is crucial for accurate personality profiling.

Emojis, while universally recognized in digital communication, can have varying interpretations even within the same language depending on the cultural context and individual user habits. In this research, the Emo-MBTI dataset was specifically focused on English-language data. However, even within the English-speaking community, the meaning and use of emojis can differ based on cultural and social factors. For example, a heart emoji might generally be interpreted as a sign of affection or love, but its use can vary from expressing platonic friendship to indicating support or solidarity depending on the context. Similarly, a thumbs-up emoji might convey agreement or approval in most contexts but could also be perceived as dismissive or sarcastic in others.

In the preprocessing of the Emo-MBTI dataset, we carefully analyzed the usage patterns of emojis within the English-language text to capture these subtle nuances. This allowed the models to understand not just the direct meaning of the emojis but also the specific emotional tone or intent that might be conveyed in different contexts. By focusing on the context in which emojis were used, the models were able to interpret the nuances of emoji-based communication more accurately. This is particularly important for personality detection, where the emotional undertones carried by emojis can provide significant insights into an individual's personality traits. For instance, frequent use of positive emojis might be associated with more extroverted or optimistic personality types, while more neutral or less frequent use of emojis might correlate with introverted or reserved traits. This careful consideration of emoji usage within the English-language Emo-MBTI dataset ensured that the models could accurately detect and interpret the emotional and personality signals conveyed by users in their digital communication.

### 3.6.3 Audio and Video Preprocessing

For datasets containing audio and video clips, specialized preprocessing steps were implemented to extract meaningful features:

- **Audio Preprocessing:**
  - *Extraction of Audio Data:* Audio was extracted from video clips and saved in WAV format.
  - *Feature Extraction:* Several features were derived from the audio data, including Mel Spectrograms, Chromograms, and Spectral Bandwidth. These features represent various aspects of the audio signal, which are critical for understanding the speaker's personality traits.
  - *Tools Used:* The librosa library was employed for audio processing, enabling the conversion of raw audio into analyzable formats.
- **Video Preprocessing:**
  - *Frame Extraction:* Three frames were extracted from each video at consistent intervals to effectively represent the visual content.



- *Image Processing*: Frames were processed using OpenCV to reduce noise and enhance clarity. This step included resizing and cropping the images to ensure consistency and quality.
- *Integration with Text Data*: The processed video frames were integrated with the corresponding textual data to provide a multimodal analysis of personality traits.

That is, the pre-Processing of Video and Textual Modalities in Chalearn dataset involved the following additional steps: For video processing, we extracted three frames from each video at consistent intervals to effectively represent the video content. Alongside the frames, audio was extracted in WAV format. This audio data was then transformed into several formats suitable for Convolutional Neural Network (CNN) [103] analysis:

- i. Mel Spectrogram Represents the power spectral density of sound.
- ii. Chromogram: Indicates the intensity of each of the 12 different pitch classes.
- iii. Spectral Bandwidth: Describes the width of the spectral energy distribution.

We employed various libraries such as OpenCV for image operations (noise reduction, resizing, cropping) and librosa for audio processing to facilitate these transformations and processing tasks. The transcripts were directly used from pickle files with dataset annotations. After filtering out punctuations, special characters, irrelevant tokens, and URLs, the data was tokenized using BertTokenizer and DebertaTokenizer. This step involved adding special tokens like [CLS] (start of a sequence), [SEP] (end of a sequence), and [PAD] (padding) to facilitate model training and inference. These tokens are essential for guiding the model about the structure of the text and improving the handling of sequence data.

#### 3.6.4 Language-Specific Preprocessing for Hindi and Bangla

Given the unique linguistic characteristics of Hindi and Bangla, additional preprocessing steps were necessary to handle these datasets effectively:

- **Text Cleaning for Hindi and Bangla**: The text was cleaned using language-specific NLP tools, such as iNLTK for Hindi, which provided functions tailored to the intricacies of Indian languages.
- **Handling Special Characters**: Special characters unique to Hindi and Bangla were carefully managed to preserve the linguistic integrity of the datasets.
- **Normalization**: Words were normalized to ensure consistency across different forms and dialects within the languages.

### **3.7. Chapter Summary**

Chapter 3 details the diverse spectrum of datasets used in our study, each contributing uniquely to the comprehensive analysis of personality traits. These datasets, spanning English, Hindi, and Bangla languages, reflect the multifaceted nature of personality expression across various cultures and communication modes. The meticulous selection, validation, and cleaning processes ensure the data's integrity, enhancing the robustness and applicability of our research findings across multiple domains. This chapter sets the foundation for exploring how personality traits are expressed and perceived in different contexts, underlining the innovative and cross-disciplinary nature of our study.

## CHAPTER 4

### **RO1: DEVELOPING AND VALIDATING CROSS-LINGUISTIC PERSONALITY DETECTION MODELS USING NATIVE LANGUAGE DATASETS AND EMOJI INTEGRATION**

Research Objective 1 (RO1) focuses on developing and validating cross-linguistic personality detection models. This objective addresses the challenge of creating AI models that can accurately identify personality traits across various languages, enhancing their adaptability and effectiveness in a global context.

A critical aspect of this research was ensuring that the models could accurately process and analyze data in different languages without compromising the cultural and linguistic nuances inherent in each language. The deep learning models discussed in this chapter, such as HindiPersonalityNet and ByaktitbaNet, were specifically designed to work with native language datasets in Hindi and Bangla, respectively. These datasets were provided to the models directly, without any translation into English. This approach allowed the models to leverage the full richness of the linguistic features present in the original texts. For instance, the HindiPersonalityNet model utilized BioWordVec embeddings, while the ByaktitbaNet model employed BERT embeddings tailored for the Bangla language. These embeddings were crucial in capturing the contextual nuances specific to each language, thereby enhancing the models' ability to accurately detect personality traits.

Importantly, the decision not to translate the contents of the Hindi and Bangla datasets into English was made to preserve the cultural and linguistic integrity of the data. By working with the native language data, the models were able to maintain the original context and meaning, resulting in more precise and culturally relevant personality detection.

In addition to processing text, the models also incorporated emojis as part of the input data. Emojis were treated as unique tokens within the text sequences, allowing the models to understand their contextual significance. This is particularly important in personality detection, as emojis often convey subtle emotional tones and nuances that are essential for accurately interpreting user intent and sentiment. Both the transformer-based models and classical models like CNNs and LSTMs could encode these emojis, integrating them into the broader analysis of personality traits.

Below is an outline of the models used in this objective. Each model in RO1 represents a significant advancement in computational linguistics, capable of processing complex datasets with a high degree of precision. These models are not only vital for the academic understanding of personality across languages but also have practical implications for various applications, such as personalized content recommendations and psycholinguistic research.

#### 4.1. Classical Deep Learning Models

Deep learning models like CNN, GRU, LSTM, and BiLSTM have become increasingly popular in the field of natural language processing (NLP) due to their ability to handle large amounts of data and learn complex relationships between words and phrases. Some commonly used deep learning models [38, 39] for NLP include:

- **Convolutional Neural Networks (CNNs):** CNNs have found a niche in Natural Language Processing (NLP) by adapting their image processing prowess to sequential data like text. In this realm, CNNs excel at tasks like text classification and sentiment analysis. Text data is first transformed into word embeddings, such as Word2Vec or GloVe [104], capturing word relationships. The CNN's convolutional layer then slides over these embeddings, using filters to identify patterns within n-grams of words. The resulting feature maps encapsulate diverse linguistic nuances. Subsequent max pooling reduces dimensionality while preserving key insights. Flattening and fully connected layers further distill information, culminating in a task-specific output layer. Though excellent for local features, CNNs might not grasp long-range dependencies as well as recurrent models.
- **Gated Recurrent Unit (GRU):** GRU [105], a specialized form of recurrent neural networks, has gained prominence in Natural Language Processing (NLP) for its adeptness in managing sequential data and capturing extended contextual dependencies in text. Addressing the vanishing gradient challenge, GRUs retain valuable information over sequences, making them adept for tasks like text generation, sentiment analysis, and machine translation. With update and reset gates, GRUs regulate past information and selective retention, enabling them to comprehend intricate language structures and relationships. While newer Transformer-based models with self-attention mechanisms dominate certain NLP realms, GRUs maintain significance in scenarios where efficiency and interpretability are paramount, offering a pragmatic balance between performance and complexity.
- **Long Short-Term Memory (LSTM):** LSTM [106][107] is a specialized recurrent neural network architecture that has become a crucial tool in Natural Language Processing (NLP) for its ability to process sequences and capture intricate relationships in text. LSTMs were designed to address the vanishing gradient problem in traditional RNNs, and they incorporate memory cells with gates that regulate information flow - input, output, and forget gates. These gates allow LSTMs to maintain important context across long sequences, making them highly effective for language modeling, machine translation, and speech recognition tasks in NLP. By considering the entire history of word sequences, LSTMs excel in understanding contextual nuances and uncovering complex patterns. While newer Transformer-based

models with parallel processing and attention mechanisms have gained prominence in NLP, LSTMs continue to be relevant in situations where the understanding of sequential information remains essential.

- **Bidirectional Long Short-Term Memory (BiLSTM):** BiLSTM [107] has emerged as a potent tool in Natural Language Processing (NLP) by enhancing the understanding of sequential data and capturing contextual nuances in text. Building upon the LSTM architecture, BiLSTM processes input data in both forward and backward directions, allowing it to capture dependencies from both past and future words. This bidirectional approach is particularly advantageous for tasks such as named entity recognition and sentiment analysis, where context from surrounding words is vital. By combining information from both directions, BiLSTM excels in grasping complex linguistic relationships. However, BiLSTM's computational complexity can be higher compared to traditional LSTMs due to the bidirectional nature. Despite this, its capability to exploit complete contextual information makes BiLSTM a valuable asset in NLP, complementing other advanced architectures like Transformer-based models.
- **Attention:** A mechanism that can be added to various neural network architectures, including RNNs and transformers, to allow the model to selectively focus on specific parts of the input text. Attention has been shown to improve the performance of many NLP tasks, including machine translation and sentiment analysis.

## 4.2. Word Embeddings

With the advent of deep learning, embedding layers became vital in NLP due to their ability to transform words into continuous vectors, capturing semantic relationships. Unlike traditional methods, these embeddings harnessed distributional semantics, learning from context, and allowing transfer learning. Embedding layers bridged the gap between raw text and neural networks, enabling nuanced language understanding and driving significant progress in NLP performance. There are two main categories of word embeddings [108][109] used in natural language processing (NLP): static and dynamic embeddings. Static embeddings, such as GloVe [104], Crawl [110], Wikipedia [111], BioWordVec [112], GoogleNews [108], and PubMed [113], are pre-trained on a large corpus and remain fixed throughout the NLP task. They provide a fast and efficient approach to encoding words with semantic and syntactic information.

In contrast, dynamic embeddings, such as ELMo [114], fastText [115], and BERT, are generated during the training of a neural network for a specific NLP task and evolve over time. They can capture more context-dependent information, including the latest language trends, and are optimized for the specific NLP task at hand. However, they are computationally more expensive and require more training data than static embeddings.

Static word embeddings are pre-trained numerical representations of words in a fixed-dimensional vector space. Created through unsupervised learning on large text datasets, they capture semantic relationships based on co-occurrence patterns. These embeddings facilitate tasks like sentiment analysis, translation, and classification by quantifying word meanings in a continuous vector space. Despite lacking context adaptability, they offer transferability, enhancing downstream NLP tasks.

- **Common Crawl embeddings** are word vectors pre-trained on the vast and diverse textual content of the Common Crawl project. This project involves regularly scraping and archiving web pages from across the internet. The embeddings are generated by analysing the co-occurrence patterns of words within this extensive dataset. These patterns allow the embeddings to capture semantic relationships and contextual information. Common Crawl embeddings provide a way to transfer the knowledge embedded in the wide range of web content to various natural language processing tasks, making them particularly useful for tasks that involve understanding and processing text from different sources on the internet.
- **FastText** is a word embedding model developed by Facebook's AI Research. It introduces sub-word information by breaking words into character n-grams. Word vectors are created by averaging the vectors of these sub-word units. The training objective is similar to skip-gram, predicting nearby words based on sub-word units. FastText handles out-of-vocabulary words by summing sub-word vectors. It's effective for morphologically rich languages, compound words, and rare words. While it enriches word representations, its use of sub-word units can increase model size and might not capture all word-level nuances.
- **GloVe**, or Global Vectors for Word Representation, is a word embedding model that combines global and local word co-occurrence statistics to generate word vectors. It constructs a co-occurrence matrix that reflects how often words appear together in a given context window. The model then factorizes this matrix to produce word vectors, optimizing a global objective function that captures the relationships between words. Unlike methods like Word2Vec, GloVe doesn't rely solely on local context, making it able to capture both syntactic and semantic relationships. These embeddings have been pre-trained on large corpora and can be directly used in various NLP tasks due to their semantic-rich nature.
- The **GoogleNews** word embeddings are high-dimensional vectors created by training a Word2Vec model on a large corpus of news articles collected from the Google News service. These embeddings capture semantic relationships between words and are commonly used for a variety of natural language processing tasks. Due to the size and quality of the training data,

GoogleNews embeddings offer rich semantic information and are especially useful for tasks that require general language understanding. They can be used to compute word similarities, analogies, and as features for different NLP applications.

- **PubMed embeddings** characterize specialized word vector representations meticulously crafted to cater specifically to the intricacies of biomedical and life sciences literature. Derived from the vast expanse of the PubMed repository, these embeddings exhibit a refined capacity to encapsulate domain-specific terminologies, thereby facilitating nuanced investigations within the intricate landscape of the biomedical domain. The nuanced semantic intricacies inherent in scientific discourse find meticulous representation within PubMed embeddings, rendering them pivotal in augmenting a spectrum of tasks, including medical text mining, disease classification, and precise biomedical information retrieval.

#### 4.3. Transformer-Based models

The transition from classical deep learning models to transformer-based models marks a fundamental shift in Natural Language Processing (NLP). Traditional models like RNNs and CNNs struggled with contextual nuances, but transformers introduced self-attention mechanisms that revolutionized context understanding. Models such as BERT, GPT, and T5 emerged, with BERT capturing bidirectional context, GPT excelling in text generation, and T5 framing tasks uniformly. XLNet [116] and RoBERTa [117] further refined training strategies, and ELECTRA [118] introduced efficient pre-training. Despite resource-intensive training, transformers' pre-trained embeddings have elevated NLP through transfer learning, significantly improving benchmark performance and propelling the field into an era of enhanced language comprehension.

- **BERT (Bidirectional Encoder Representations from Transformers)** is a groundbreaking transformer-based model for natural language processing (NLP). It employs a two-step pre-training process: Masked Language Model (MLM) and Next Sentence Prediction (NSP). During MLM, BERT randomly masks words in input sentences and predicts them using surrounding context, enabling bidirectional understanding. In NSP, BERT predicts whether two sentences follow each other in a document, enhancing context awareness. BERT's architecture includes multiple layers of self-attention mechanisms, allowing it to capture contextual dependencies across words in both directions. It produces contextualized word embeddings that are fine-tuned for downstream tasks like classification, question-answering, and more. BERT's impact lies in its ability to understand complex context, enabling state-of-the-art performance across a wide range of NLP tasks by leveraging large-scale, diverse pre-training data.



- **Decoding-enhanced BERT with Disentangled Attention (DeBERTa)** [119] is an advanced transformer-based model tailored for natural language processing (NLP). It introduces disentangled attention mechanisms that disentangle different types of dependencies, enhancing contextual understanding. DeBERTa further integrates decoding techniques during both pre-training and fine-tuning stages, improving generation capabilities. The model's architecture combines bidirectional and unidirectional training, mitigating the shortcomings of traditional BERT models. DeBERTa achieves state-of-the-art performance across NLP tasks by leveraging disentangled attention, multi-granular tasks, and enhanced bidirectional and unidirectional information flows. Its holistic approach to context modeling and decoding makes it a potent tool for diverse NLP applications.
- **Efficiently Learning an Encoder that Classifies Token Replacements Accurately (ELECTRA)** is a cutting-edge model designed for natural language processing (NLP). It introduces a unique pre-training task that involves replacing tokens within a sentence and training the model to differentiate between original and replaced tokens. This approach enhances pre-training efficiency by focusing on the replaced tokens. ELECTRA employs a generator-discriminator framework, where a generator creates replacements, and a discriminator judges their authenticity. This adversarial training refines the model's ability to discern meaningful tokens. By efficiently utilizing training data, ELECTRA improves upon traditional masked language model (MLM) pre-training. It has shown promising results on various NLP benchmarks, benefiting tasks like text classification, named entity recognition, and more, while also being computationally efficient.
- **Metas Open Pre-Trained Transformer Language Models (Meta OPT)** is a framework that extends pre-trained transformer-based models for natural language processing (NLP). It emphasizes the importance of adaptability and customization of these models to various tasks and domains. Meta OPT introduces a meta-learning approach, where the model is trained on a diverse range of tasks during pre-training. This prepares the model to rapidly adapt and fine-tune on new tasks with minimal data. By leveraging a shared parameter space across tasks, Meta OPT aims to achieve improved performance and efficient utilization of training resources. This framework advances the concept of transfer learning in NLP by providing a versatile tool for building domain-specific models that can quickly adapt to new challenges.
- **XLNet** is an advanced transformer-based model for natural language processing (NLP) that combines bidirectional and autoregressive training. It introduces a novel permutation-based training approach, enabling the model to capture both contextual information from surrounding words and the

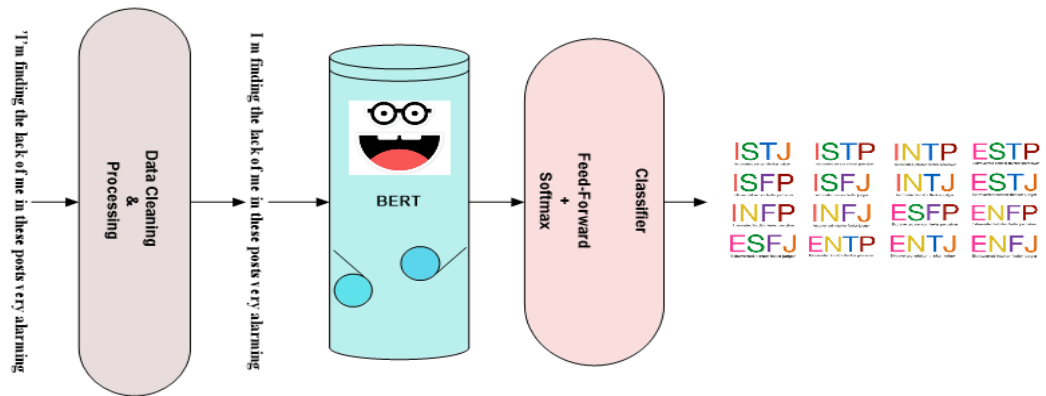


predictive nature of autoregressive models. This results in improved handling of bidirectional context and fine-grained relationships between words in text. XLNet's architecture enhances its understanding of context and relationships, leading to state-of-the-art performance on various NLP tasks such as text classification, sentiment analysis, and language generation.

#### 4.4. PersonalityBERT (English)

This part of research puts forward a BERT (Bidirectional Representation for Transformers) based model for recognizing apparent personality traits from textual modality. BERT includes a transformer with multiple attention mechanisms to provide context of words in a text and a classification layer to the transformer output to detect the personality from textual data. It is one of the most popular neural architectures used for a wide variety of NLP tasks and fine-tuning BERT allows to build a robust classification model to predict categories. The proposed PersonalityBERT is textual modality-specific deep neural model that fine-tunes a pre-trained BERT for the personality classification task.

Leveraging the BERT architecture, the PersonalityBERT model is engineered to analyze the Kaggle\_MBTI dataset, categorizing English language data into 16 MBTI personality types. The architecture model is shown in the figure 4.1 below.

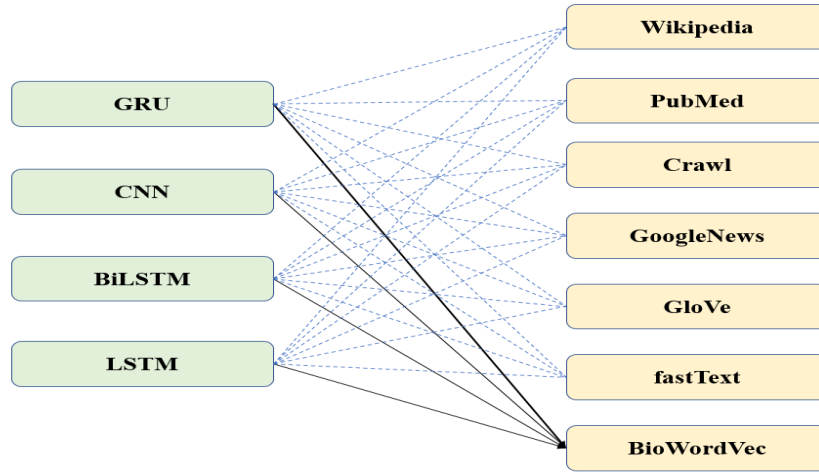


**Fig. 4.1.** Architecture of PersonalityBERT model

#### 4.5. HindiPersonalityNet (Hindi)

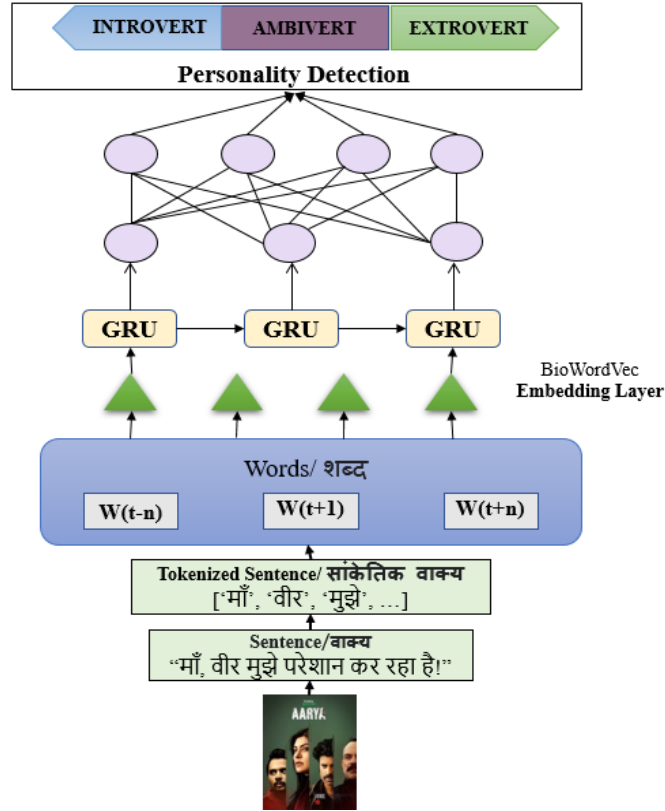
Tailored for the Hindi language, the HindiPersonalityNet employs sophisticated BioWordVec embedding with a GRU to interpret the Shaksiyat dataset, which includes conversational data from the series "AARYA." This model differentiates between introvert, extrovert, and ambivert personality expressions.

In this research, we utilized GRU and performed a comparative analysis involving LSTM, BiLSTM and CNN. Fig. 4.2 showcases the diverse combinations of embeddings and deep learning models employed in this study.



**Fig.4.2.** Deep learning and embedding employed

The model proposed HindiPersonalityNet model uses BioWordVec embedding to train GRU (Gated Recurrent Unit). BioWordVec embeddings are trained on a large corpus of biomedical literature, including scientific texts, articles, and abstracts, which enables them to learn the specific language used in the biomedical domain. This makes them effective at capturing the nuances of language use in specific domains, which can be beneficial for personality detection in conversational data. GRU is a type of recurrent neural network (RNN) architecture that is capable of processing sequential data. It is designed to overcome some of the limitations of traditional RNNs, such as the vanishing gradient problem, which can make it difficult to train deep models. GRUs use a gating mechanism to selectively update the hidden state at each time step, allowing them to capture long-term dependencies in the input sequence. To build the NLP model, we first initialize the GRU with random weights and then train it on the BioWordVec embeddings. The model takes in a sequence of word embeddings as input, which are processed by the GRU layer. The output of the GRU layer is then fed into a fully connected layer for the classification tasks. Fig.4.3 depicts the system architecture of the proposed model.



**Fig.4.3.** The HindiPersonalityNet Model

#### 4.6. ByaktitbaNet (Bangla)

In 1947, psychologist Eysneck [120] added another category as *ambivert personality* for individuals showing traits of introverted personality in some situations and in others, they behave as an extrovert personality type. Fig.4.4 depicts the characteristics of each of these personality type.

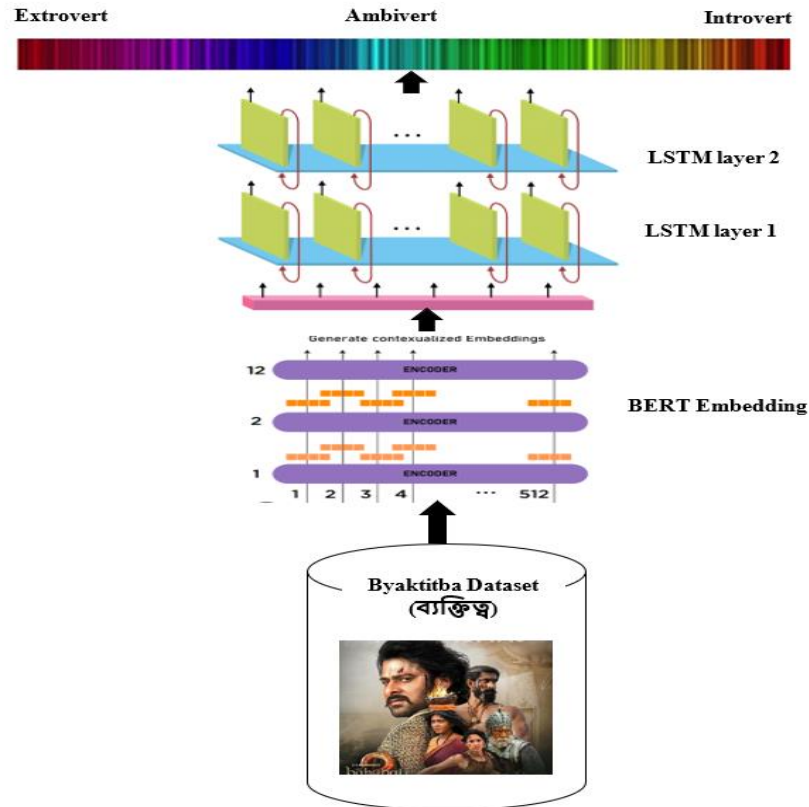


Extrovert	Ambivert	Introvert
<ul style="list-style-type: none"> <li>• People and social situations excite and energize them.</li> <li>• Usually initiate and engage in conversations.</li> <li>• Can talk about anything with anyone.</li> <li>• Don't mind others paying full attention to them.</li> <li>• Meeting new people doesn't faze them.</li> </ul>	<ul style="list-style-type: none"> <li>• Often wonder whether they need alone time or external stimulation.</li> <li>• Could be quiet during the entire conversation, but also share what they are passionate about.</li> <li>• If in the right context, don't mind attention, but often they prefer standing at the side-lines.</li> <li>• Are fine with talking to new people, but it's better to do it with their friends.</li> </ul>	<ul style="list-style-type: none"> <li>• Alone time is the way to recharge.</li> <li>• Use their eyes and ears more than their mouths.</li> <li>• Don't like small talks.</li> <li>• Prefer standing away from the spotlight.</li> <li>• Quite uncomfortable to meet new friends.</li> </ul>

**Fig.4.4.** The personality continuum spectrum

As individual tendency towards the outer world (extroversion) or the inner world (introversion) can have a major influence on career choice, relationships and overall lifestyle, this research puts forward a personality detection model for assorting dialogues into three personality categories {ambivert, extrovert, introvert} in conversational data using dialogues from the Bengali version of the Indian epic action film, Baahubali 2: The Conclusion (বাহুবলি ২: দ্যা কনক্লিউশান). This research is a preliminary work that considers the personality continuum spectrum to categorize personalities based on the dimension of attitude as extrovert, introvert and ambivert. We propose a BERT-LSTM model, ByaktitbaNet (ব্যক্তিত্ব, meaning personality in Bengali). BERT (Bidirectional Encoder Representations from Transformers) optimally captures the contextual features of the dialogues and the structural and semantic details of the Bengali language. The generated feature vectors are then passed through stacked LSTM (Long short-term memory) layers for processing them further and exploring the hidden inter-dependencies amongst the words of the dialogues. The stacked LSTM layers also classify the dialogues into their respective categories as ambivert, extrovert or introvert.

The ByaktitbaNet model is specifically developed for the Bangla language, applying BERT embeddings combined with LSTM to process the Byaktitba dataset derived from "Baahubali 2" transcripts. The model excels in classifying text into introvert, extrovert, and ambivert categories, reflecting personality expression in Bangla. The architecture of this model is shown in the fig. 4.5.



**Fig.4.5.** ByaktitbaNet Model

#### 4.7. Transformer-based Models for *Personality\_Quotes* (English)

In this part of our research centres on classical deep learning models, specifically Convolutional Neural Networks (CNN), Long Short-Term Memory (LSTM), Gated Recurrent Unit (GRU), and Bidirectional Long Short-Term Memory (BiLSTM). Additionally, we delve into static embeddings: Common Crawl, GloVe, FastText, PubMed, and GoogleNews. These embeddings complement the aforementioned deep learning models in predicting personalities for personality quotes. Further, we discuss five most recent transformer-based models, namely ELECTRA, BERT, DeBERTa, Meta OPT and XLNet.

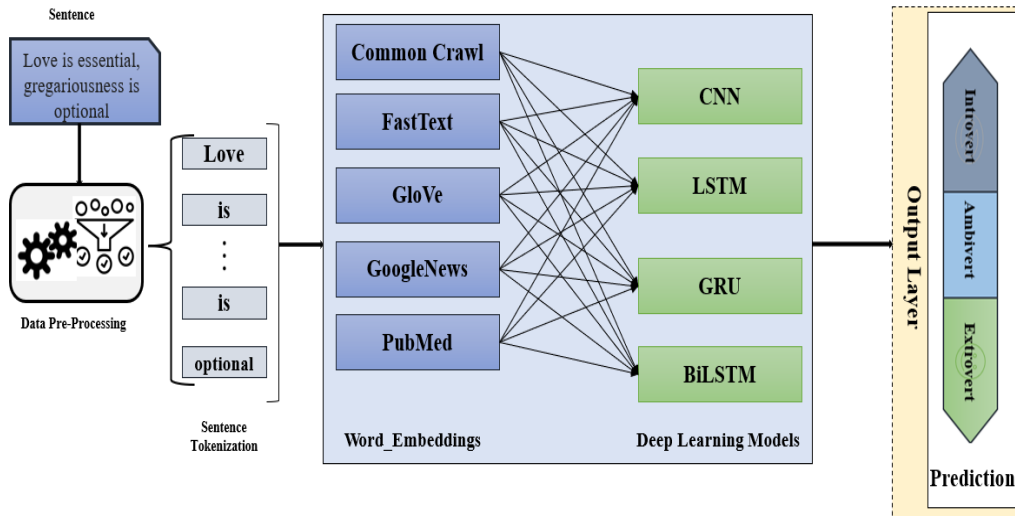
The use of deep learning and transformer models on a small dataset like "*Personality\_Quotes*" is justified for several reasons:

- **Complex Relationships:** Deep learning and transformer models are designed to capture intricate relationships and patterns in data, even when the dataset is relatively small. They excel at learning representations from data with multiple layers of abstraction, which is beneficial when dealing with complex concepts like personality traits.
- **Feature Extraction:** Deep learning and transformer models can automatically extract relevant features from raw text data. This is particularly

advantageous when dealing with unstructured text, as seen in personality quotes, where manually crafting features might be challenging or insufficient.

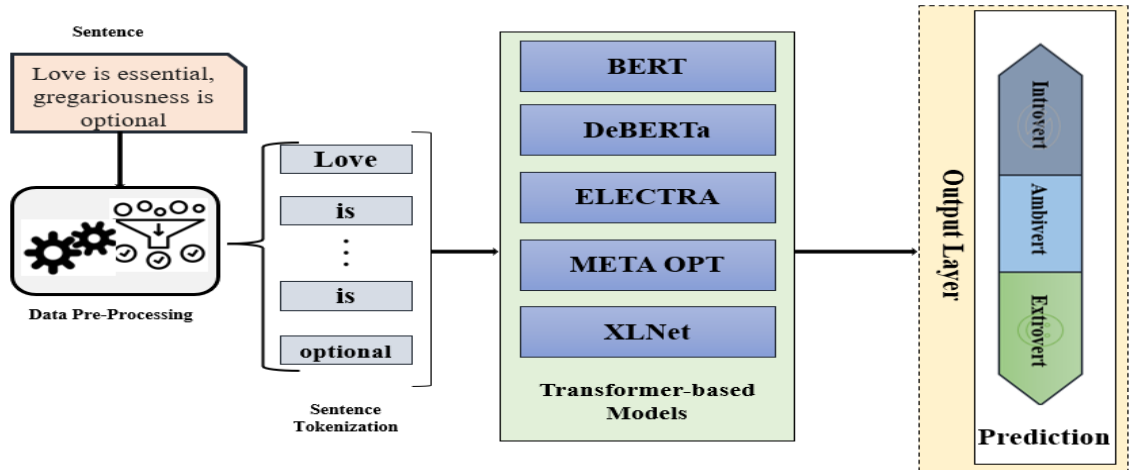
- **Transfer Learning:** Transformer-based models like BERT, ELECTRA, and DeBERTa are pre-trained on large text corpora, which enables them to generalize well to smaller datasets. The pre-trained knowledge can be fine-tuned on the specific task of personality detection, making them effective even with limited data.
- **Performance Improvement:** Deep learning models have shown remarkable performance improvements in various natural language processing tasks, even on small datasets. This suggests that their ability to capture intricate linguistic nuances can lead to meaningful insights even with limited training examples.
- **Potential for Generalization:** While small datasets may raise concerns about overfitting, the transfer learning and regularization mechanisms present in deep learning and transformer models can mitigate this issue. This allows the models to learn from the data while also incorporating general language understanding.
- **Rich Representations:** Transformers are known for generating rich contextualized embeddings that can capture the meaning of words in the context of the entire sentence. This is highly beneficial when dealing with personality detection, where context plays a crucial role in understanding traits.
- **Framework for Future Data:** Training deep learning and transformer models on small datasets can lay the groundwork for future improvements as more data becomes available. The models can be retrained or fine-tuned with larger datasets to further enhance their accuracy.
- **Comparative Analysis:** The use of these advanced models on a small dataset can provide a comparative analysis of their effectiveness. This can shed light on their potential benefits, limitations, and generalization capabilities in scenarios with limited data.

Our approach harnesses the strengths of classical deep learning models paired with static word embeddings like Common Crawl and GloVe. This combination is leveraged to predict personalities within the "*Personality\_Quotes*" dataset, showcasing how these traditional models can still provide valuable insights in the age of transformers. Fig.4.6 diagrammatically summarizes the pairings between classical deep learning models and static embeddings used.



**Fig.4.6.** Word Embeddings with Deep Learning Models

To complement the classical architectures, we employ the latest transformer-based models, including BERT, DeBERTa, Meta OPT, and XLNet. These models have been fine-tuned on our "Personality\_Quotes" dataset to capture a broader and more contextually rich representation of personality traits. Figure 4.7 shows the transformer-based models evaluated.



**Fig.4.7.** Personality Prediction using Transformer-based Models

#### 4.8. Chapter Summary

This chapter has comprehensively detailed the methodologies and foundational models used to achieve Research Objective 1 (RO1), focusing on the development and validation of cross-linguistic personality detection models. By employing a robust array of classical deep learning architectures such as CNNs, GRUs, LSTMs, and BiLSTMs alongside advanced transformer-based models including BERT,

DeBERTa, ELECTRA, Meta OPT, and XLNet, we have demonstrated the capability to process and analyze complex multilingual datasets with high precision. The integration of static and dynamic word embeddings has further enhanced the contextual understanding and semantic richness of these models. Each technique and model explored within this objective not only advances the field of computational linguistics but also provides substantial practical implications for applications in personalized content recommendation systems, psycholinguistic research, and beyond. This multifaceted approach underscores the potential of these models to adapt and excel in a global context, marking a significant leap forward in the realm of cross-linguistic personality detection.

This chapter introduced innovative models such as HindiPersonalityNet and ByaktitbaNet, designed to process native language datasets without translation. It also explores the integration of emojis into deep learning models for enhanced personality detection. These contributions represent significant advancements in the field, addressing the challenges of cross-linguistic and emotion-based personality analysis.



## CHAPTER 5

### **RO2: INTEGRATING PSYCHOLOGICAL THEORIES INTO DEEP LEARNING FRAMEWORKS FOR MULTI-LANGUAGE AND CULTURAL CONTEXTS**

Research Objective 2 is focused on enriching personality detection models through the integration of deep psychological theories and psychometric methods. In this segment of our research, we navigate the intricacies of psychological constructs, like parenting styles, and correlate them with MBTI personality types using data from low-resource linguistic datasets. This objective is particularly ambitious, as it aims to blend the qualitative depth of psychological assessments with the quantitative strength of advanced computational models.

In this chapter, we explore the integration of psychological theories with advanced deep learning frameworks, focusing on personality detection across different linguistic and cultural contexts. The chapter makes several key contributions: the development of a kernel-based soft-voting ensemble model evaluated on both English and Hindi datasets, the creation of a novel Decision-maker MBTI dataset for psycholinguistic profiling, and the introduction of predictive modelling approaches for predicting parenting styles within Hindi-speaking populations. These contributions significantly expand the applicability of AI-driven personality detection models in diverse environments, enhancing our understanding of how personality traits manifest across different cultures and languages.

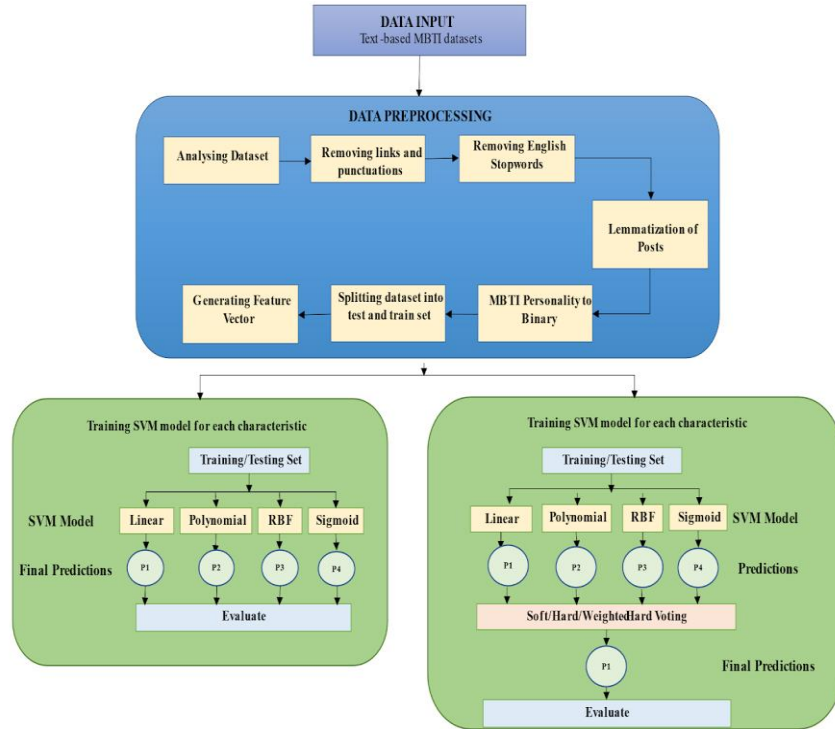
#### **5.1. Kernel-based Soft Voting Ensemble Model for Personality Detection (KBSVE-P)**

Section 5.1 introduces the Kernel-based Soft Voting Ensemble Model (KBSVE-P), which has been specifically designed to enhance personality detection across English and Hindi datasets. This section discusses the model's architecture, the innovative use of various SVM kernels, and the application of advanced voting techniques to improve prediction accuracy in diverse linguistic contexts.

The novel kernel-based soft-voting ensemble model for personality detection from natural language, KBSVE-P is built by firstly evaluating the performance of various Support Vector Machine (SVM) kernels, namely radial basis function (RBF) [121], linear, sigmoidal, and polynomial, to find the best-suited kernel for automatic personality detection in natural language text. Next, an ensemble of SVM kernels is implemented with a variety of voting techniques such as soft voting, hard voting, and weighted hard voting. The model is evaluated on the publicly available Kaggle MBTI dataset for detecting a user's personality using various performance metrics. Simultaneously, a MBTI based Hindi personality dataset, विशेष चरित्र\_MBTI (pronounced as vishesh charitr, meaning personality in Hindi) is built and the KBSVE-P model is evaluated on this dataset too.

The Count Vectorizer transformation is used to convert tokenized textual data into a token count matrix. Term Frequency-Inverse Document Frequency (TF-IDF) [122] transformation is used to convert the matrix into a normalized TF-IDF representation. It is used to analyze how much a word is relevant in the collection of sentences. Here, TF is an estimate of the frequency of a word that appears in a document, while IDF represents how important a word is in a set of documents. It is calculated based on how rarely a word appears in a set of documents.

The architectural flow of the proposed a kernel-based ensemble model for personality detection, KBSVE-P is shown in fig. 5.1. The sub-sections discuss the details of the architecture.



**Fig. 5.1.** Kernel-based Ensemble Model (KBSVE-P) evaluated on English MBTI and Hindi विशेष चरित्र\_MBTI datasets.

### 5.1.1. Support Vector Machine (SVM)

The purpose of an SVM is to classify the given data into two classes by finding an appropriate hyperplane that can separate the classes effectively. SVM uses different linear and non-linear kernels to find the best-suited hyperplane. Support Vector machines can be used in case of non-regularly distributed data and data with unknown distribution; hence it is suitable for textual data. It efficiently finds hyperplanes to separate two classes present in any distribution due to the availability of different kernels that define equations for hyperplanes. The SVM kernel used for the classification of textual data includes the Linear Kernel, Radial Basis function kernel (RBF), Sigmoid kernel, and Polynomial Kernel. In SVM, kernel implies

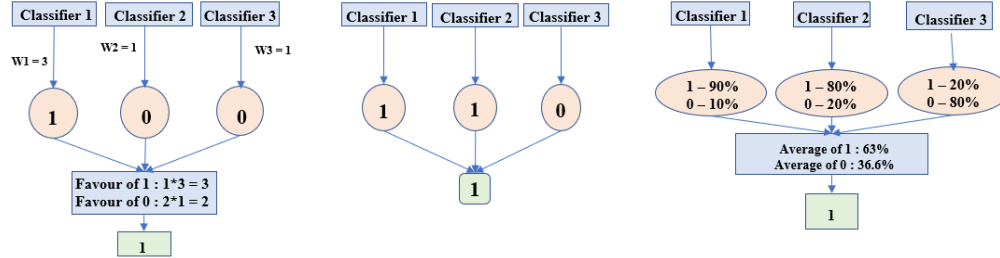
internal feature transformation. It is a way to compute the dot product of two vectors in a higher dimensional feature space. That is, the kernel method helps in calculating the dot product in a different space without even visiting the original space. The time taken to compute the standard dot product is  $O(n^2)$ , whereas the kernel takes only  $O(n)$  time to calculate the dot product. RBF is used in SVMs to help SVM to become non-linear rather than linear. A linear kernel is mainly used when we have a linearly separable set of data points. A polynomial kernel is a non-stationary kernel, which is well suited for problems where all the training data is normalized.

To find the best suited SVM kernel for classification from textual data, we train four SVM models, one model for each kernel. The Kernels mentioned above are the same kernels used for analyzing the performance of an ensemble of SVM kernels aggregated with the three voting techniques, namely soft voting, hard voting, and weighted hard voting.

### 5.1.2. Voting

Voting classifiers are ensemble methods for decision-making and are further divided into 2 categories: soft and hard voting. Hard voting entails picking the prediction with the highest number of votes, whereas soft voting entails combining the probabilities of each prediction in each model and picking the prediction with the highest total probability. This voting [123] classifier built by combining different classification models turns out to be a stronger meta-classifier that balances out the individual classifier's weakness on a particular dataset. It takes majority voting based on weights applied to the class or class probabilities and assigns a class label to a record based on the majority vote. A soft voting classifier is one in which the output class is predicted based on the average probability given to that class by different classifiers. The class with the highest average probability will be the final output of the ensemble. In soft voting, every individual classifier provides a probability value that a specific data point belongs to a particular target class. The predictions are weighted by the classifier's importance and summed up. Then the target label with the greatest sum of weighted probabilities wins the vote. Assume that the prediction probability of class A by 3 different classifiers is given as [0.30, 0.47, 0.53] and the prediction probability of class B is [0.20, 0.32, 0.40]. Therefore, the average prediction probability of class A is 0.43 and B is 0.306. Hence the output of the ensemble aggregated with soft voting will be class A. In hard voting classification, the output class is predicted based on the number of votes given to the concerned class. The class with the majority of votes is the output class. In hard voting, the predictions of each algorithm are considered with the ensemble selecting the class with the highest number of votes. For example, if three algorithms predict the color of a particular wine as white, white, and red, the ensemble will predict red. It is also called majority voting. Suppose that the predicted output of 3 different classifiers is class A, A, and B; it results in the majority predicted class 'A' as an output. Hence, the output of the ensemble aggregated with hard voting will be class A. Weighted hard voting is an extension of hard voting where different classifiers may have different weights. The class with the majority votes will be predicted as the output

class, but every classifier may have different weights associated with them, which means the output of a particular model may be more significant than others. Fig. 5.2 depicts diagrammatically representing these voting classifiers.



**Fig.5.2.** Voting Classifiers (Soft Voting, Hard Voting, and Weighted-Hard Voting)

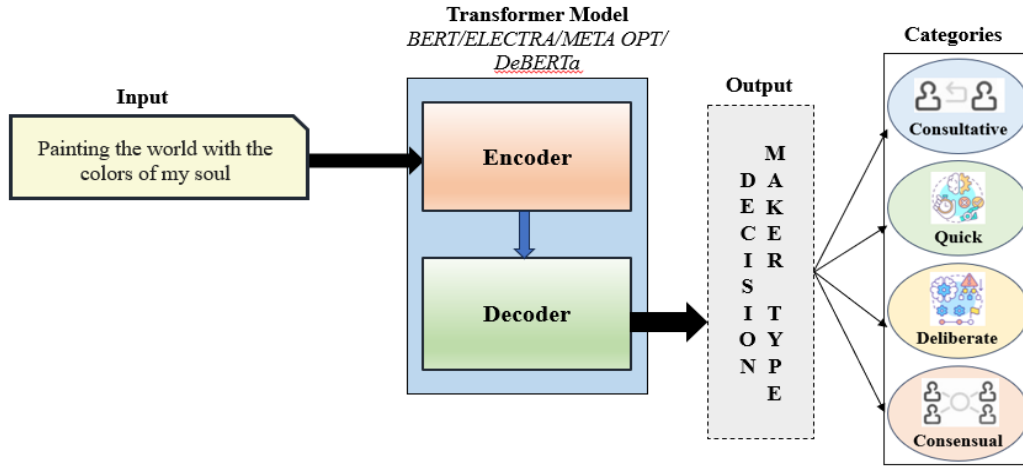
Thus, the KBSVE-P model evaluates predictions through a dual-phase approach: initially using individual SVM kernels and subsequently applying soft, hard, and weighted-hard voting techniques. This ensemble strategy not only amplifies prediction accuracy but also imparts resilience to the model, enabling it to handle the diverse and complex nature of personality data effectively. The KBSVE-P model aligns with RO2 by leveraging advanced machine learning techniques to enhance personality detection models. By integrating psychological theories (MBTI) with computational methods, this model exemplifies the interdisciplinary approach intended to combine qualitative psychological insights with quantitative data analysis, thereby enriching the overall personality detection process.

## 5.2. Transformer-Based Models for Decision-Maker-MBTI Type

Section 5.2 presents the creation and application of the Decision-maker MBTI dataset, which is re-annotated for psycholinguistic profiling. This section also details the implementation of state-of-the-art Transformer models to predict decision-making styles based on MBTI personality types, highlighting how these models can be effectively used in English-language datasets

This part of research has a theoretical component, which focuses on describing a conceptual framework for correlating personality traits to decision-maker types. Based on the identification of key decision-maker categories from literature review, a novel dataset, Decision-maker- MBTI (DM-MBTI) dataset was created to identify personality traits associated with specific decision-maker types. As discussed previously, to mitigate ethical and compliance concerns related to data privacy and protection, we repurposed publicly available text-based personality detection datasets through a careful re-annotation process. This re-annotation aligned with established literature and was also validated by a consultant psychologist. The experimental component included developing a deep learning model for automated MBTI decision-maker type identification, that harnesses the capabilities of state-of-the-art Transformer models prediction model to recognize the decision-maker type of the subject based on the user-centric psychometric NLP. Transformer

models, including BERT, DeBERTa, ELECTRA and Meta OPT are employed to identify MBTI decision-maker types based on user-generated content. (Fig.5.3).



**Fig.5.3.** MBTI Decision-Maker Type Recognition using Transformer Models evaluated on re-annotated English MBTI dataset

The following Transformer-based models were implemented and evaluated for classifying the posts into four decision maker categories:

- **BERT (Bidirectional Encoder Representations from Transformers):** BERT is a transformer-based model that employs a masked language model and next sentence prediction during pre-training. It captures bidirectional context by learning to predict missing words within sentences and understanding the relationships between consecutive sentences.
- **DeBERTa (Decoding-enhanced BERT with Disentangled Attention):** DeBERTa improves upon BERT by introducing disentangled attention mechanisms, which help the model disentangle different types of dependencies within the input text. This enhances the model's ability to capture complex contextual relationships.
- **ELECTRA (Efficiently Learning an Encoder that Classifies Token Replacements Accurately):** ELECTRA introduces a novel pre-training approach using a generator-discriminator framework. Instead of predicting masked words, it trains the model to distinguish between original and replaced tokens, making pre-training more efficient.
- **Meta OPT (Metas Open Pre-Trained Transformer Language Models):** Meta OPT is a framework that emphasizes adaptability and customization. It utilizes a meta-learning approach during pre-training, exposing the model to a

wide range of tasks, enabling it to quickly adapt and fine-tune for specific tasks with limited data.

Thus, this research conducts AI-enhanced psycholinguistic profiling for decision style prediction by utilizing MBTI and transformer models. It explores how advanced AI and psycholinguistic techniques can predict individuals' decision-making styles. By integrating the Myers-Briggs Type Indicator (MBTI) framework with state-of-the-art transformer models, this study aims to develop robust and accurate methods for profiling decision styles, offering significant applications in personalized learning, targeted marketing, and improved human-computer interaction. The work aligns with RO2 by demonstrating how modern Transformer-based models can be utilized to understand and predict decision-making styles through MBTI personality traits. This approach underscores the fusion of deep psychological theories with state-of-the-art machine learning techniques, further advancing the interdisciplinary goals of RO2.

### **5.3. Machine Learning for Predicting Parenting Styles via MBTI**

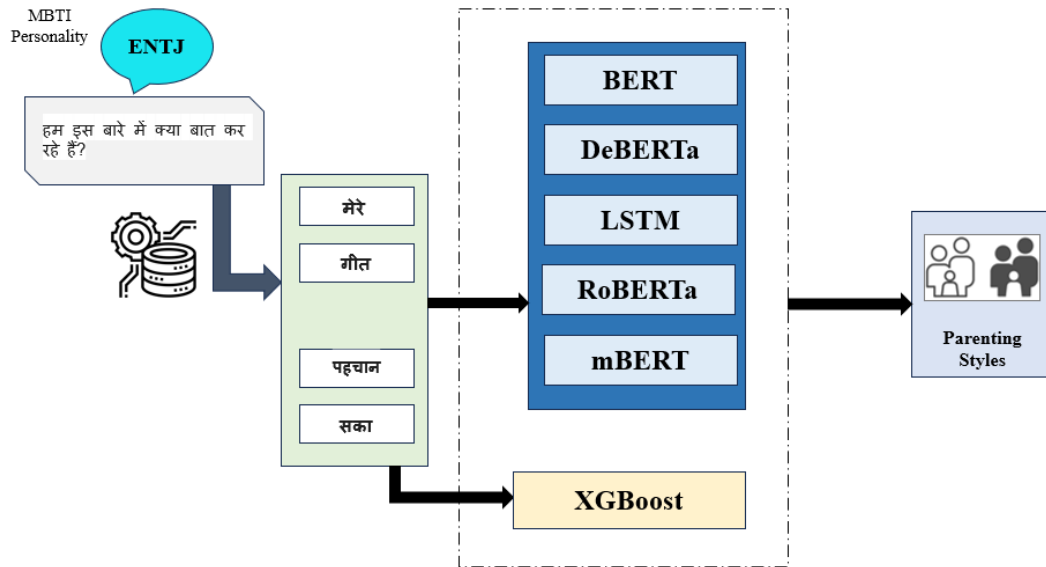
Section 5.3 explores the relationship between parenting styles and MBTI personality types within Hindi-speaking populations. This section introduces a predictive modelling approach that evaluates deep learning models to compare their performance against a tree-ensemble model, XGBoost, to predict parenting styles, thus emphasizing the role of cultural context in psychological assessment.

#### **5.3.1. Parenting Styles and MBTI Profiles**

Parenting styles, distinguished by approaches to discipline, communication, expectations, and nurturing, significantly influence child development. In the context of MBTI-derived personality types, "*Energetic Explorers*" (ENFP, ESFP, ENTP) and "*Steadfast Guardians*" (ISFJ, ISTJ, ESTJ) provide unique perspectives on parenting behaviours. Energetic Explorers are enthusiastic, spontaneous, and adventurous, fostering independence and curiosity in their children but sometimes lacking structure. Steadfast Guardians are reliable and structured, emphasizing rules and traditions but potentially being rigid and limiting their children's independence.

Other parenting styles include authoritative (balanced expectations and support), authoritarian (strict and controlling), permissive (lenient and indulgent), and uninvolved (neglectful and detached). Integrating MBTI profiles into these styles provides a nuanced view of how personality influences parenting, helping tailor approaches to individual family dynamics and cultural contexts. The fig.5.4 and 5.5 illustrate our approach:

- The fig.5.4 shows the use of different models; machine learning algorithm (XGBoost); Deep Learning model (LSTM) and transformer-based models (BERT, DeBERTa, RoBERTa, mBERT) to predict personality traits.



**Fig.5.4.** Workflow of Parvarish Dataset with different models for Personality-Parenting Prediction

This study bridges the gap between personality assessment and cultural psychology, employing a predictive modelling approach to predict parenting styles in Hindi-speaking populations. By integrating psychometric analysis with cultural psychology, the research underscores the importance of cultural context in psychological research and highlights the potential of psychometric tools in understanding human behaviour across diverse cultural landscapes. This approach not only enhances our understanding of personality development but also allows for the cultural adaptation of psychological theories, ensuring their relevance and applicability across diverse populations. Predictive modelling based on MBTI profiles provides nuanced perspectives on family dynamics, aiding in the development of tailored parenting interventions that respect cultural values and practices. Thus, this segment aligns with RO2 by integrating psychometric analysis and cultural psychology with machine learning techniques to predict parenting styles. This convergence of methodologies exemplifies the interdisciplinary approach intended to blend psychological theories with advanced computational models, thereby enriching the understanding and application of personality psychology in diverse cultural contexts.

#### 5.4. Chapter Summary

This chapter underscores the interdisciplinary endeavour of integrating psychological theories with machine learning techniques. The findings enhance our understanding of personality detection and broaden the applicability of AI in psychological research and practical applications. By leveraging advanced computational models alongside deep psychological insights, the study illustrates the potential of AI to transform personality psychology. Future research directions include refining these models and



exploring their potential in various domains such as personalized learning, targeted marketing, and improved human-computer interaction. The integration of qualitative psychological assessments with quantitative AI techniques represents a significant advancement, promising innovative applications in both academic research and real-world settings.

This chapter has made several significant contributions to the field of computational personality detection. By integrating psychological theories with deep learning models, we have enhanced personality detection across different languages and cultural contexts. The development of the KBSVE-P model and its application to both English and Hindi datasets highlight the model's versatility and effectiveness. Additionally, the creation of the Decision-maker MBTI dataset and the introduction of predictive modelling for predicting parenting styles in Hindi-speaking populations offer new insights into the intersection of psychology and AI. These advancements not only broaden the scope of personality detection but also provide practical applications in personalized learning, targeted marketing, and cultural psychology.



## CHAPTER 6

### RO3: PROFILING EMOTIONAL DISPOSITIONS AND EVALUATING REAL-WORLD APPLICATIONS

Research Objective 3 addresses the intricate task of profiling emotional dispositions such as optimism and pessimism, and it extends to appraising the real-world applicability of our models in domain such as online personalized employment, and education. We also scrutinize the capability of transformer models to discern MBTI personality types through the lens of emoji usage in digital communication.

This chapter makes significant contributions to the field of computational personality detection by introducing several novel methodologies and approaches:

- **Chakra System for Emotional Profiling:** This chapter introduces a pioneering 'Chakra' system that maps the interaction between personality traits and emotional attitudes, providing a novel framework for understanding human psychological dynamics.
- **Advanced Transformer-Based Models:** The chapter details the development and fine-tuning of advanced transformer models, including mBERT, XLM-RoBERTa, IndicBERT, and a novel stacked mDeBERTa, specifically designed to profile emotional dispositions in Hindi textual data. These models represent a significant advancement in the ability to analyze complex emotional and personality data across different languages.
- **Emoji-Based Personality Prediction:** The research further explores the use of emojis in digital communication as a novel data source for predicting MBTI personality traits, integrating these visual cues with textual data to enhance psychological profiling in online interactions.

These contributions highlight the chapter's innovative approach and its significance in advancing both the theoretical and practical applications of personality detection in digital environments.

#### 6.1. Transformer-Based Models for Emotional Profiling

This section of research undertakes a comprehensive analysis of optimistic and pessimistic tendencies present within Hindi textual data, employing transformer-based models. The research represents a pioneering effort to define and establish an interaction between the personality and attitude chakras within the realm of human psychology. Introducing an innovative "Chakra" system to illustrate complex interrelationships within human psychology, this work aligns the Myers-Briggs Type Indicator (MBTI) personality traits with optimistic and pessimistic attitudes,

enriching our understanding of emotional projection in text. The study employs meticulously fine-tuned transformer models specifically mBERT [124], XLM-RoBERTa, IndicBERT, mDeBERTa [125] and a novel stacked mDeBERTa trained on the novel Hindi dataset ‘मनोभाव’ (pronounced as Manobhav). The following sub-sections discuss the overall human psychological make-up in terms of "Behaviour," "emotion," "attitude," and "personality" finally introducing the proposed human psychology “Chakra” system.

#### 6.1.1. Overall Human Psychological Make-up

"Behaviour," "emotion," "attitude," and "personality" are interconnected concepts that collectively contribute to understanding human psychology and individual differences. While they are distinct, they influence and shape each other in complex ways.

- **Behaviour:** Behaviour refers to the observable actions, reactions, and responses of an individual in various situations. It encompasses both verbal and nonverbal actions. Behaviour is influenced by a combination of internal factors (such as thoughts, emotions, and attitudes) and external factors (such as environmental stimuli and social context).
- **Emotion:** Emotion refers to a complex psychological state that involves feelings, physiological changes, and cognitive responses. Emotions are often triggered by internal or external stimuli and play a significant role in guiding behaviour. Emotions are subjective experiences that can range from joy and happiness to sadness, anger, fear, and more.
- **Attitude:** Attitude refers to a person's overall evaluation, belief, or opinion about a particular object, person, group, idea, or situation. Attitudes can be positive, negative, or neutral and are shaped by an individual's experiences, beliefs, values, and emotions. Attitudes influence how a person perceives and responds to the world around them.
- **Personality:** Personality is a relatively stable and enduring pattern of thoughts, feelings, behaviours, and characteristics that distinguish one individual from another. It encompasses a wide range of traits, including cognitive, emotional, and behavioural tendencies. Personality traits contribute to an individual's consistency in behaviour and emotional responses across different situations and over time.

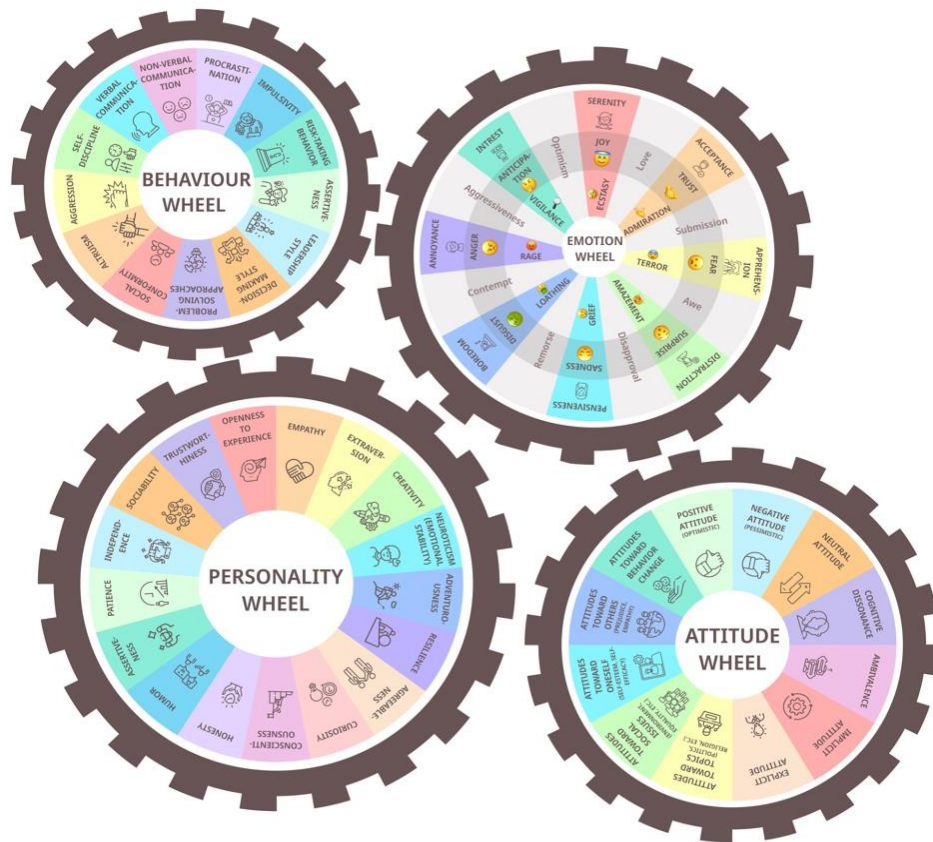
In a cohesive fusion, these elements intertwine to craft a nuanced tapestry of human existence. Within this narrative, emotions delicately influence attitudes, attitudes intricately mould personality, and personality seamlessly directs behaviour. The relationships among these concepts can be described as follows:

- **Behaviour and Emotion:** Emotions can drive behaviour. For example, feeling excited may lead to energetic and enthusiastic behaviour, while feeling anxious might result in cautious or avoidant behaviour. Similarly, behaviours can influence emotions. Engaging in positive behaviours may lead to increased feelings of happiness.
- **Behaviour and Attitude:** Attitudes can shape behaviour. If someone has a positive attitude toward exercise, they are more likely to engage in physical activity. On the other hand, behaviours can also impact attitudes. Repeatedly engaging in a behaviour may lead to attitude change through cognitive dissonance or reinforcement.
- **Emotion and Attitude:** Emotions can influence attitudes. For example, experiencing positive emotions in a certain context may lead to the development of a positive attitude toward that context. Attitudes can also influence emotions. Having a negative attitude toward a task may result in feelings of frustration or stress when engaging in that task.
- **Personality and behaviour/Emotion/Attitude:** Personality traits play a role in shaping behaviour, emotions, and attitudes. For instance, an extraverted individual might engage in more social behaviours and experience positive emotions in social settings. Certain personality traits may also predispose individuals to specific attitudes, such as a more open-minded attitude in individuals with high openness to experience.

In summary, behaviour, emotion, attitude, and personality are intertwined elements of human psychology. They interact and influence each other in intricate ways, contributing to the complexity of individual human experiences and expressions.

### 6.1.2. Human Psychology Chakra System

The word "chakra" means "wheel" or "disk" in Sanskrit. A chakra model is a well-established conceptual framework stemming from ancient Indian traditions, like Yoga and Hinduism, depicting the body's energy centres as spinning wheels or vortexes. It represents the body's energy centres, often depicted as spinning wheels or vortexes of energy that correspond to specific physical, mental, emotional, and spiritual aspects of human existence. Now visualizing this intricate interplay within human psychology as a chakra system with four wheels, each representing a crucial facet: personality, emotion, attitude, and behaviour reveals a captivating analogy. Like spinning wheels of energy, these harmoniously interact, contributing to the holistic balance of our inner emotional and mental landscape as shown in fig.6.1.



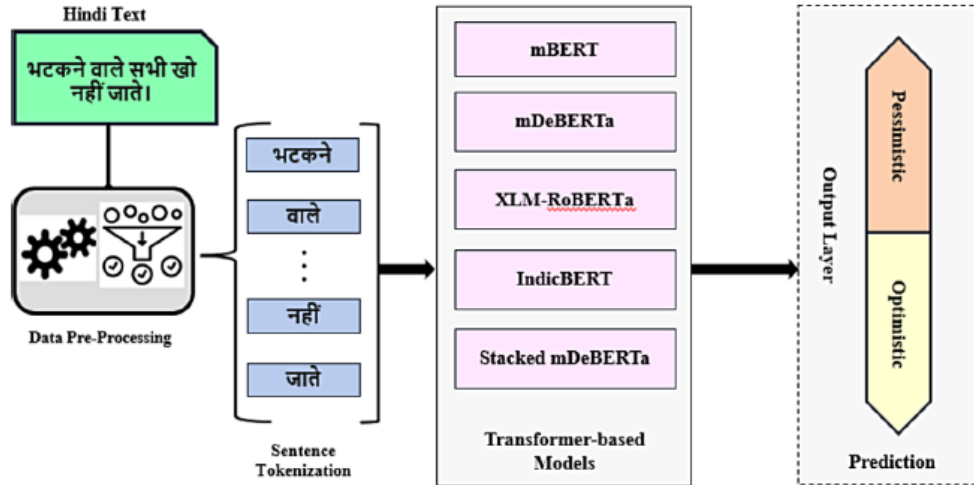
**Fig.6.1.** The Human Psychology Chakra System

- **Personality Wheel:** This foundational wheel embodies the essence of who we are, housing our enduring traits, tendencies, and core characteristics. It rotates with the rhythms of our unique nature, influencing how we perceive and respond to the world around us.
- **Emotion Wheel:** Enveloping the personality wheel, this vibrant sphere radiates with the hues of our inner feelings. It spins with the dance of joy, sorrow, excitement, and more, infusing our experiences with a vivid spectrum of emotional shades. Psychologist Robert Plutchik developed one of the most popular emotion wheels, known as the Plutchik wheel [126].
- **Attitude Wheel:** Nurtured by the emotion wheel, this wheel of perspective spins, shaping how we interpret and engage with life's tapestry. It revolves with optimism, pessimism, or neutrality, casting a distinct filter through which we view our surroundings.
- **Behaviour Wheel:** The outermost wheel, powered by the synergy of personality, emotion, and attitude, propels our outward actions and

interactions. It rotates with precision, orchestrating our responses and movements, translating inner dynamics into tangible expressions.

In this harmonious chakra system, the personality wheel serves as the core, while the emotion, attitude, and behaviour wheels revolve around it. Just as the balanced flow of energy in a chakra system promotes well-being, the fluid interaction of these four wheels nurtures a holistic understanding of human psychology, guiding our journey through self-discovery and personal growth.

For the fine-grained task of emotional disposition analysis, we harness a selection of transformer-based models as depicted in the provided diagram. These models, including mBERT, mDeBERTa, XLM-RoBERTa, IndicBERT, and a stacked configuration of mDeBERTa, are strategically employed to process Hindi text. After rigorous data pre-processing and sentence tokenization, these models feed into an output layer that classifies the input as either pessimistic or optimistic, thereby providing a deep understanding of the underlying emotional attitudes. Fig.6.2 shows the transformer-based models evaluated.



**Fig.6.2.** Emotional Attitude prediction using transformer-based models

- mBERT (Multilingual Bidirectional Encoder Representations from Transformers) Model:** Multilingual BERT (mBERT), developed by Google AI, is a transformative variant of the original BERT model that is designed to tackle the challenges of multilingual natural language processing (NLP). Unlike the standard BERT, which was mainly trained on English text, mBERT is pretrained on a diverse and expansive collection of text data in multiple languages. This multilingual pretraining equips mBERT with the ability to capture language-agnostic patterns, syntax, and context that transcend individual languages. The BERT multilingual base model (uncased) model is pre-trained on 102 languages and some of them are: Georgian, German, Gujarati, Hindi, Hungarian and etc. mBERT's multilingual nature empowers it with cross-lingual transfer capabilities. It

learns to recognize similarities in linguistic structures and relationships that hold true across different languages. Moreover, mBERT offers an intriguing advantage in zero-shot and few-shot learning scenarios. This means that the model can be fine-tuned with just a small amount of labelled data from a new language, even if it isn't explicitly trained on that language.

- **mDeBERTa (Multilingual Decoding-enhanced BERT with disentangled attention) Model:** mDeBERTa is multilingual version of DeBERTa which use the same structure as DeBERTa and is trained with CC100 (Common Crawl 100) multilingual data. It is an advanced transformer-based language model that represents a significant evolution of the original mBERT architecture which is mainly based on BERT architecture. It is developed to address some limitations of traditional transformer models and further enhance the performance of natural language processing tasks. The mDeBERTa V3 base model have 12 layers and a hidden size of 768. It has 86M backbone parameters with a vocabulary containing 250K tokens which introduces 190M parameters in the Embedding layer. This model was trained using the 2.5T CC100 data.
- **XLM-RoBERTa (Multilingual RoBERTa) Model:** XLM-RoBERTa, developed by Facebook AI, combines BERT and RoBERTa (A Robustly Optimized BERT Pretraining Approach) models for advanced multilingual language understanding. It's trained to comprehend text across languages, using robust pre-training techniques inspired by RoBERTa. The model excels in cross-lingual tasks, learning from diverse languages. It improves contextual understanding through masked language modelling and uses byte-pair encoding for versatile tokenization. XLM-RoBERTa adapts well to specific languages via fine-tuning and employs unsupervised tasks like translation and cross-lingual modelling. Overall, it's a powerful and language-agnostic model for multilingual natural language processing.
- **IndicBERT Model:** IndicBERT is a specialized variant of the BERT model that has been tailored and optimized for Indic languages, a group of languages spoken in the Indian subcontinent. Similar to other BERT-based models, IndicBERT aims to capture rich contextual information from text data, making it well-suited for a wide range of natural language processing (NLP) tasks. The distinct feature of IndicBERT lies in its fine-tuning and customization to better handle the linguistic nuances, structures, and characteristics specific to Indic languages. This involves training the model on large and diverse datasets from Indic languages, enabling it to learn patterns and representations that are particularly relevant to these languages.

IndicBERT is a multilingual ALBERT [127] model pretrained exclusively on 12 major Indian languages. It is pre-trained on a novel monolingual corpus of around 9



billion tokens and subsequently evaluated on a set of diverse tasks. IndicBERT has much fewer parameters than other multilingual models (mBERT, XLM-R etc.) while it also achieves a performance on-par or better than these models. The 12 languages covered by IndicBERT are: Assamese, Bengali, English, Gujarati, Hindi, Kannada, Malayalam, Marathi, Oriya, Punjabi, Tamil, Telugu.

### 6.1.3. Proposed Stacked mDeBERTa Model

Stacking two mDeBERTa models involves employing two instances of the mDeBERTa architecture in a sequential manner. Each instance can be thought of as a separate model, and by stacking them, they form a composite architecture with specific benefits and applications. Some of the potential details and implications of stacking two mDeBERTa models are:

- **Increased Representational Power:** Stacking two models allows for a more complex and powerful representation of input data. The second model can capture higher-order relationships and nuances in the outputs of the first model, leading to more expressive representations.
- **Feature Hierarchies:** The first mDeBERTa model might capture lower-level features and structures, while the second model can build upon those features to capture more abstract and high-level patterns in the data.
- **Fine-tuning and Specialization:** The first model can serve as a general feature extractor, and the second model can be fine-tuned for specific tasks or domains. This approach can lead to better task-specific performance.
- **Error Correction:** Stacking models can help mitigate errors or inconsistencies that might occur in the predictions of a single model. The second model might learn to correct errors made by the first model, enhancing overall accuracy.
- **Ensemble Effect:** Stacking can be seen as a form of model ensemble, where the combined predictions of two models can lead to improved overall performance by reducing bias and variance.

It's important to note that while stacking models can offer advantages, it increases computational requirements, as each model adds to the processing load. At the same time, hyperparameters of both models need to be tuned carefully to achieve optimal performance. This includes architectural choices, learning rates, and regularization parameters.

This research seamlessly aligns with Research Objective 3 (RO3) by advancing our understanding of emotional dispositions, such as optimism and pessimism, through the use of sophisticated transformer-based models like mBERT, XLM-RoBERTa, IndicBERT, mDeBERTa, and a novel stacked mDeBERTa. By

introducing the innovative "Chakra" system to map the interplay between Myers-Briggs Type Indicator (MBTI) traits and these emotional attitudes, the study enriches our comprehension of how personality and emotion are expressed in text. This alignment is further demonstrated by the models' application in diverse domains such as online shopping, personalized employment, and education, showing how these insights can be translated into real-world applications. The research not only deepens the theoretical understanding of the complex relationships among behaviour, emotion, attitude, and personality but also showcases practical applications, thus perfectly encapsulating the essence of RO3 in bridging advanced computational approaches with practical, real-world utility.

## **6.2. Fusion Model: XceptionResNet with BERT for Chalearn (Multimodal) Dataset**

Using the First Impressions V2 dataset, we developed a fusion model that integrates XceptionResNet [128] with BERT to adeptly extract and analyze complex personality traits from video data. This hybrid model not only demonstrates superior performance in multimodal personality analysis but also shows how such methodologies can enhance applications in smart city environments. The model processes video and transcript data through a systematic approach, utilizing Convolutional Neural Networks (CNNs) for frames and audio, and transformer models for textual content, ensuring comprehensive analysis. By merging these modalities, the fusion model effectively learns from visual, auditory, and textual cues, paving the way for advanced applications in fields like personalized employment and behavioural prediction, aligning perfectly with RO3's goal of applying emotional profiling in diverse real-world settings.

In this research, the Chalearn First Impression Dataset is utilized to provide a comprehensive understanding of video and transcript processing using systematic approaches. This dataset, a significant resource in personality analysis, consists of 10,000 videos from over 3,000 HD YouTube videos, each labelled with the Big Five personality traits: Extraversion, Agreeableness, Conscientiousness, Neuroticism, and Openness. Each video, averaging 15 seconds in length, is accompanied by detailed transcripts that average 43 words per video, with extensive annotations reflecting various sentiments and personality dimensions.

To enhance the analysis, the videos were processed to extract three key frames at consistent intervals and their audio in WAV format. These frames and audio data were then transformed into formats such as Mel Spectrogram, Chromagram, and Spectral bandwidth, suitable for Convolutional Neural Networks (CNNs). This transformation involved meticulous pre-processing steps including noise reduction, resizing, and cropping, using tools like OpenCV, PIL, and librosa, to ensure the integrity and clarity of the data.

For textual analysis, the transcripts were directly used from the datasets pickle files, pre-processed using tokenizers to filter out punctuations and irrelevant tokens. This led to the creation of a CSV file combining text transcripts with the corresponding Big Five personality trait labels, which were tokenized using



BertTokenizer and DebertaTokenizer to enhance model compatibility. This structured approach not only allowed for a deeper exploration of the multimodal data but also demonstrated how different data modalities can be integrated to predict personality traits effectively, aligning with the broader objectives of enhancing behavioural prediction and personality analysis in real-world applications.

### 6.2.1. Deep Neural Architectures

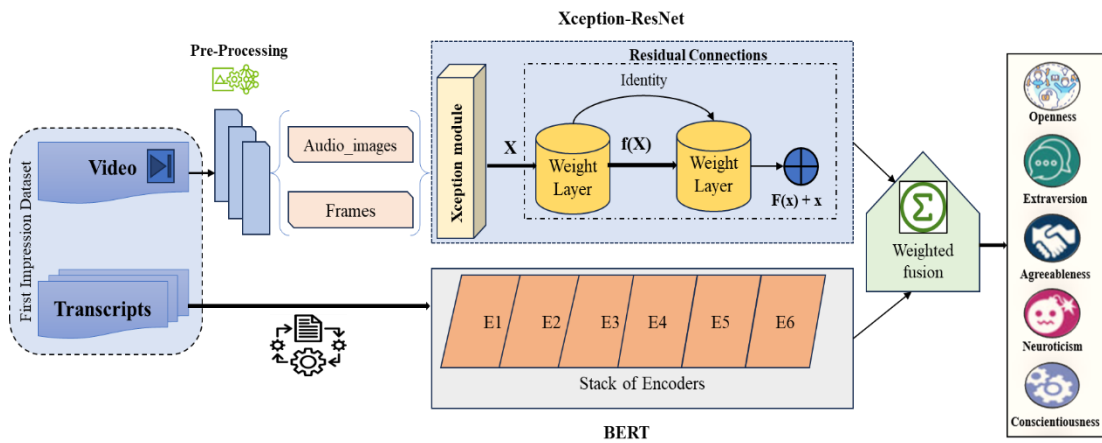
When it came to model application, a varied yet targeted approach was adopted. The video modality, which encompasses both frames and the transformed audio images, was analysed using models such as InceptionV3 [129], EfficientResNet [130], and XceptionResNet [128]. In contrast, the textual content derived from transcripts was put through transformer models like BERT and DeBERTa. Our recent paper demonstrates the applicability of transformer-based models in understanding complex personality traits from textual data. Recognizing the potency of a combined approach, fusion models were developed, seamlessly integrating features from video-focused models with the textual capabilities of BERT, ensuring a comprehensive multimodal data analysis. With our processed video frames and audio feature images in hand, the next crucial step is data fusion. We fuse the processed video frames and audio feature images into a single input tensor. This integration of modalities allows our model to learn from both visual and auditory cues simultaneously.

- **InceptionV3:** InceptionV3 uses inception modules with various filter sizes and factorized 7x7 convolutions to capture features at different scales, improving computational efficiency and accuracy in vision tasks.
- **EfficientResNet:** EfficientResNet combines the scaling capabilities of EfficientNet and the deep network training benefits of ResNet, incorporating skip connections for balanced accuracy and computational cost.
- **XceptionResNet:** XceptionResNet leverages the Xception architecture's depth wise separable convolutions and ResNet's residual connections to capture complex patterns more efficiently.
- **BERT:** BERT revolutionized NLP by training on a large corpus using a masked language modelling objective, capturing deep contextual relationships between words for a wide range of NLP tasks.
- **DeBERTa:** DeBERTa enhances BERT with disentangled self-attention mechanisms, improving contextual comprehension and capturing long-range dependencies in sequential data modelling.

### 6.2.2. Multimodal Fusion Architectures

This section presents a comprehensive overview of the fusion models employed in personality trait prediction. In the preceding section, the deep neural architectures

described are now employed in an integrated manner, incorporating the technique of weighted fusion. Weighted fusion is a valuable approach applied to amalgamate data from diverse modalities within the realm of multimodal data analysis. This technique grants the capability to allocate differing levels of significance to the contributions of each modality, contingent upon their pertinence and effectiveness in a given task. By doing so, it enhances the overall performance and resilience of multimodal systems, ensuring that the most pertinent and impactful information from various sources is effectively harnessed for comprehensive analysis and prediction. The first architecture, Inception V3 + BERT, combines InceptionV3 and BERT models to simultaneously process textual, image, and audio data. The video modality is processed using InceptionV3, while the text modality undergoes transformation by the BERT transformer. The concatenated outputs from these models pass through a series of densely connected layers to predict the Big Five traits. Similarly, the EfficientResNet + BERT architecture combines EfficientResNet and BERT for multimodal data processing. EfficientResNet handles image inputs, while BERT processes text inputs. These outputs are fused using a weighted fusion technique and subsequently processed through dense layers to predict personality traits. The final fusion architecture, XceptionResNet + BERT, utilizes Xception blocks with residual connections for image inputs and pretrained BERT for text data, ultimately employing weighted fusion and dense layers to predict personality traits. Fig.6.3 illustrates the XceptionResNet + BERT multimodal fusion architecture.



**Fig.6.3.** XceptionResNet + BERT Structure

Incorporating findings from the First Impression dataset used in this study, we can leverage advanced AI models to decode subtle cues in candidates' expressions, speech patterns, and behaviours that signal underlying personality traits. This method enables a nuanced understanding of how a candidate's personality aligns with the cultural and functional demands of a job role. By applying a model trained on this dataset, recruiters can systematically evaluate and match candidates to positions where they are most likely to thrive, thereby optimizing job satisfaction and organizational productivity. This approach not only streamlines the recruitment

process but also ensures a deeper harmony between an individual's innate characteristics and their professional responsibilities, embodying the essence of personality-driven employment.

The fusion model combining XceptionResNet with BERT, developed using the First Impressions V2 dataset, exemplifies the alignment of sophisticated multimodal analysis with Research Objective 3 (RO3). This model adeptly integrates video and textual data, applying Convolutional Neural Networks (CNNs) and transformer models to extract and interpret complex personality traits from diverse modalities. The systematic approach, which processes frames, audio, and transcripts, results in a comprehensive understanding of personality traits, enhancing applications in smart city environments and fields such as personalized employment and behavioural prediction.

By employing deep neural architectures like InceptionV3, EfficientResNet, XceptionResNet, BERT, and DeBERTa, the model benefits from a robust, varied analysis approach. These architectures allow for the effective integration of visual, auditory, and textual cues, facilitating advanced personality trait predictions. The fusion of these data through weighted fusion techniques optimizes the use of each modality, ensuring that the most pertinent and impactful information is harnessed for a thorough analysis.

This fusion model's success in processing and analyzing multimodal data underlines its potential to revolutionize personality analysis in real-world applications. By aligning with RO3's goals, this model demonstrates how leveraging advanced AI techniques can lead to deeper insights into human behaviour and personality, thereby enhancing the effectiveness of employment matching and contributing to the broader field of behavioural analysis. The ability to decode subtle behavioural cues and align them with job roles shows the model's potential to improve job satisfaction and organizational productivity, embodying the essence of personality-driven employment. This approach not only streamlines recruitment processes but also ensures a harmonious match between individuals' innate characteristics and their professional responsibilities.

### **6.3. Emo-MBTI: Mapping Emojis to MBTI Personalities**

User-generated content on social media platforms has become a valuable resource for personality detection, as it encompasses more than just text. Beyond written posts and comments, user-generated content includes images, videos, likes, shares, and other interactions, providing a rich source of data for inferring individuals' personality traits. Thus, a multi-modal approach to personality detection, incorporating textual, visual, and emotive elements, allows a more comprehensive understanding of users' personalities and preferences, facilitating personalized experiences and targeted interventions in various digital contexts. More specifically, emojis, beyond their colourful appearance, are integral to digital communication, offering visual cues that enhance the expression of emotions and attitudes alongside text. This symbiosis provides a deeper understanding of individuals' digital self-expression. For example, in a text message conversation, a person might use a series

of heart emojis 🧡💖💕 to convey affection and excitement about a forthcoming event. These emojis complement the text, adding layers of emotional nuance and conveying the sender's enthusiasm more vividly than words alone. Similarly, a laughing emoji 😂 might accompany a joke or humorous comment, amplifying the sense of amusement and creating a more engaging interaction.

As emojis play a crucial role in enriching digital communication by conveying emotions and attitudes in a concise and visually appealing manner, enhancing the overall expressive capacity of written text. analyzing emoji usage to infer personality types based on MBTI principles merges the realms of digital communication and psychological typology. This integration offers a novel approach to understanding the psychological underpinnings of individuals through their digital expressions.

Building on this foundation, this research aims to refine predictive models that adeptly correlate specific emojis with MBTI traits to discern personality types from textual conversations. For instance, frequent use of emotive emojis like hearts or smileys may suggest a preference for Feeling over Thinking, while the use of more ordered or structured symbols like checkmarks or clocks might indicate a Judging rather than Perceiving inclination. This methodology leverages the inherent expressiveness of emojis to tap into subtle psychological cues that may be less apparent in plain text. The utilization of emoji-based personality detection has practical applications across various domains. In social media, understanding user personality can enhance content personalization and ad targeting, improving user engagement. In educational technologies, it can help tailor learning experiences to suit different personality-driven learning styles, potentially increasing effectiveness and satisfaction. Furthermore, in professional settings, this approach could assist in team formation and leadership strategies by providing insights into team members' preferred communication styles and decision-making processes.

### **6.3.1. EmoMBTI-Net: The Emoji-Based MBTI Personality Prediction Model**

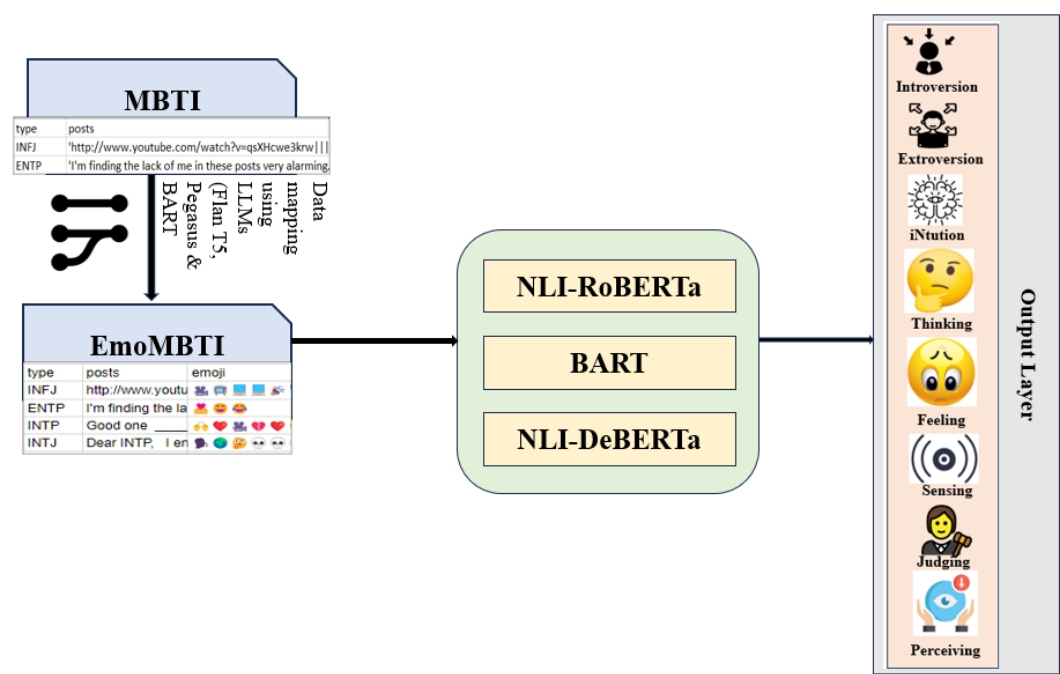
In this study, our goal is to predict MBTI personality types based on the usage of emojis in textual communications. As introduced earlier, the integration of emojis within natural language processing frameworks offers a novel approach to understanding complex human behaviours and personality traits. The rationale behind this methodology lies in the richness of emojis as expressive tools that convey emotional and contextual nuances beyond the capacity of plain text. Emojis encapsulate a spectrum of emotions and sentiments, making them valuable in the psychological profiling inherent to personality assessment.

The principal challenge we faced was selecting appropriate models equipped to process emojis as fundamental components of language. Traditional models like BERT, DistilBert[131], and ALBERT, while robust in general language tasks, do not inherently support emoji interpretation within their tokenization processes. This limitation necessitates the exploration of alternative models that natively incorporate emoji understanding to ensure the accuracy of personality predictions.

- **NLI-RoBERTa:** NLI-RoBERTa, an extension of the RoBERTa (Robustly optimized BERT approach) model, is designed for Natural Language Inference (NLI) tasks. The model employs a transformer architecture known for its efficiency in grasping contextual relationships within text. NLI-RoBERTa is pre-trained on extensive text datasets using a self-supervised learning paradigm, enabling it to capture complex semantic subtleties essential for discerning between entailment, contradiction, and neutrality in texts. Such capabilities make NLI-RoBERTa highly effective for sequence classification, including MBTI personality prediction, where understanding the layered meanings conveyed by emojis is crucial. The strength of NLI-RoBERTa lies in its transformer-based design, incorporating multiple layers of self-attention mechanisms and feedforward neural networks. This architecture is adept at learning contextual embeddings from vast amounts of data, crucial for emoji-based text interpretation. By eliminating the next sentence prediction task and adopting dynamic masking, the model enhances its focus on relevant textual segments, crucial for accurate personality assessment.
- **BART:** BART [132] is versatile in handling both sequence-to-sequence and sequence classification challenges. Its architecture combines a bidirectional transformer encoder with an auto-regressive decoder, forming a denoising autoencoder that excels in reconstructing noisy inputs to predict original sequences. This innovative training approach enables BART to identify and retain critical textual features, including emojis, making it suitable for detailed personality analysis. The bidirectional encoder ensures comprehensive understanding of context, capturing dependencies in text from both directions, which is vital for interpreting the emotional and contextual layers conveyed by emojis. The auto-regressive decoder enhances the model's ability to generate coherent outputs that are contextually aligned with the input, thereby supporting sophisticated sequence classification tasks like MBTI personality prediction.
- **NLI-DeBERTa:** NLI-DeBERTa refines the DeBERTa model with a focus on Natural Language Inference. Its standout feature is the disentangled attention mechanism, which optimizes focus on different segments of input text independently, crucial for parsing complex emoji-laden communications. This ability to focus selectively on relevant parts of the text facilitates a deeper understanding of long-range dependencies and intricate semantic relationships, essential for accurate personality profiling. The specialized architecture of NLI-DeBERTa, with its nuanced attention to detail in context modelling, makes it an exceptional choice for tasks requiring an advanced grasp of subtle textual interactions. Its efficacy in MBTI classification is enhanced by its capacity to disentangle attention weights during pre-training,

thus improving the model's overall ability to interpret and analyze emoji-based communication in the context of personality assessment.

The following fig.6.4 illustrates a framework for analyzing personality types using advanced NLP models, specifically designed to process and interpret digital communication related to MBTI personality traits through emoji usage. It showcases a system that inputs data from various MBTI and emoji-based posts (EmoMBTI), which are then processed by NLP models like NLI-RoBERTa, BART, and NLI-DeBERTa. These models analyze text and emojis to interpret underlying personality traits, as indicated by the MBTI categories on the right side of the image, including Introversion, Extroversion, Intuition, and others. The diagram effectively communicates how modern AI tools can be utilized to derive insights about personality from online interactions, suggesting a methodical approach to psychological profiling in digital environments.



**Fig.6.4.** The MBTIEmoNet Model

This research aims to harness emojis as a novel data source for deducing personality traits based on the Myers-Briggs Type Indicator (MBTI), enhancing our grasp of digital behavior and its psychological connotations. The study is structured around several key objectives: This involves curating and refining a unique dataset that integrates emoji usage with textual data from existing MBTI datasets and additional posts scraped from the r/mbti subreddit on Reddit. The goal is to map emojis to specific MBTI personality types, thereby enriching the tools available for personality analysis. The project employs state-of-the-art language models such as



FlanT5, BART, and Pegasus to probe and elucidate the contextual relationships between text and emojis. This phase is crucial for laying the foundation necessary for precise personality type predictions. The research utilizes finely tuned transformer models like RoBERTa, DeBERTa, and BART to predict MBTI personality types from emoji usage. By leveraging the deep understanding of emoji contexts, this approach aims to significantly boost the accuracy of personality profiling, demonstrating the potential of emojis in psychological assessments in digital communications.

Emo-MBTI: Mapping Emojis to MBTI Personalities, aligns perfectly with Research Objective 3 (RO3) by exploring the role of emojis in digital communication to infer MBTI personality types, thus advancing our understanding of digital behaviour and its psychological underpinnings. The use of user-generated content, which includes not only text but also emojis, images, and videos, provides a rich dataset for analyzing personality traits in a multimodal context. This research leverages emojis as a significant, expressive tool within digital dialogues, recognizing that these visual cues enhance textual communication by adding emotional and attitudinal depth that plain text might lack.

By developing the Emo-MBTI dataset, which integrates emoji usage with textual data, this study taps into a novel resource for personality analysis that could transform how personality traits are inferred online. Additionally, implementing and evaluating advanced language learning models like FlanT5, BART, and Pegasus to map emoji-text contextual relationships and fine-tuning transformers to predict personality types showcases a sophisticated methodological approach. These steps ensure a comprehensive analysis capable of capturing the nuanced interplay between textual content and visual emojis, reflecting the study's alignment with RO3's aim to apply emotional profiling in varied real-world settings. This integration not only enriches the theoretical landscape of personality psychology but also enhances practical applications such as personalized digital marketing, educational technologies, and workplace optimization, making significant strides towards understanding and leveraging personality dynamics in digital interactions.

#### **6.4. Chapter Summary**

This chapter focuses on RO3, which aims to profile emotional dispositions such as optimism and pessimism and evaluate the real-world applications of various models in domains like online shopping, personalized employment, and education. This chapter particularly highlights the use of transformer-based models to analyze Hindi textual data for emotional profiling and introduces an innovative "Chakra" system to better understand the interplay of personality traits with emotional attitudes through a structured analysis of human psychological make-up behaviour, emotion, attitude, and personality. The chapter also outlines the integration of emoji with textual data to map emojis to specific MBTI personality types, enhancing personality analysis capabilities. It discusses the implementation and evaluation of advanced language models like FlanT5, BART, and Pegasus to understand the contextual relationship between text and emojis, aiming to refine predictive models for accurate personality

type predictions. In a significant section on deep neural architectures, the text elaborates on using various models like InceptionV3, EfficientResNet, and XceptionResNet for analyzing video modalities, and transformer models like BERT and DeBERTa for textual content. This multimodal approach underscores the integration of visual, auditory, and textual data through advanced fusion techniques, which improve the understanding and prediction of personality traits.

Furthermore, the chapter addresses the practical application of these methodologies in enhancing user interaction in digital platforms through personalized experiences, showcasing how sophisticated AI models can be employed to understand and predict human behaviour and personality traits effectively. The strategic use of emojis in digital communication, as discussed, not only enriches user interaction but also provides deeper insights into personality types, demonstrating a novel approach to psychological profiling in digital environments.



## CHAPTER 7

### FINDINGS, DISCUSSION, AND LIMITATIONS

In this chapter, we discuss the detailed results of our research, examining how different models perform across our distinct research objectives. We take a model-wise look within each objective to understand the effectiveness of our approaches in personality detection. As we move through the chapter, we assess each model's contributions, providing a clear view of their capabilities.

#### 7.1. Performance Metrics used in Research

Following are the metrics used in our research.

- **Accuracy:** It is defined as the ratio of correct predictions (TP and TN) to the total number of samples (all entries of the confusion matrix summed up). Equation (7.1) defines the accuracy, which is as follows:

$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN} \quad (7.1)$$

- **Precision:** It measures how many of the samples predicted as positive are actually positive. Thus, it is defined as the ratio of true positive (TP) values to the summation of all predicted positive values (TP and FP). Equation (7.2) specifies the precision, which is as follows:

$$\text{Precision} = \frac{TP}{TP + FP} \quad (7.2)$$

- **Recall:** It measures how many of the positive samples are captured by the positive predictions. Therefore, it is defined as the ratio of true positive (TP) values to the summation of all actual positive values (TP and FN). Equation (7.3) characterizes the recall, which is as follows:

$$\text{Recall} = \frac{TP}{TP + FN} \quad (7.3)$$

- **F1-Score:** It is also termed as F-score and is an important metric for measuring a model's performance. It is defined as the harmonic mean of precision and recall. Equation (7.4) defines the F1-Score, which is as follows:

$$\text{F1-Score} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} \quad (7.4)$$

- **ROC-AUC:** ROC-AUC is the Area Under Receiver Operating Characteristic Curve that summarizes the performance of a classifier over all possible thresholds. The ROC curve plots the true positive rate also known as recall on the y-axis against the false positive rate on the x-axis.
- **MCC:** MCC stands for Matthews Correlation Coefficient, which measures the correlation coefficient between the predicted values and actual values. Its values range between -1 and +1. A coefficient of +1 represents a perfect

prediction, 0 an average random prediction, and -1 in an inverse prediction. Equation (7.5) specifies the MCC, which is as follows:

$$\text{MCC} = \frac{TP * TN - FP * FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (7.5)$$

where; TP: True Positive, TN: True Negative, FP: False Positive, and FN: False Negative

- **Mean Squared Error (MSE):** It quantifies the average squared difference between predicted and actual values in a regression model, providing a measure of the model's prediction accuracy.
- **Root Mean Squared Error (RMSE):** It is the square root of the Mean Squared Error and is a better measure for the accuracy of a model.

### 7.1.1 Model Training and Validation

To ensure the accuracy and reliability of the performance evaluations presented in this chapter, it is essential to establish that each model was properly trained and validated. This section details the training processes that were implemented for each model discussed. For all models, the following steps were rigorously applied:

- **Dataset Splitting:** The datasets were divided into training and testing sets using an 80-20 split. In certain cases, a 10-fold cross-validation technique was applied to further validate the model's performance and ensure its generalizability across different subsets of the data.
- **Hyperparameter Tuning:** Key hyperparameters such as learning rate, batch size, and the number of epochs were carefully tuned to optimize model performance. Grid search and random search methods were utilized to find the optimal hyperparameter settings.
- **Regularization Techniques:** To prevent overfitting, which can lead to misleadingly high performance on the training data but poor generalization to new data, regularization techniques such as dropout and early stopping were employed. These techniques ensure that the model does not learn noise from the training data and performs well on unseen data.
- **Model Validation:** The performance of each model was validated using the test set, which was not seen by the model during training. The use of validation data provided an unbiased evaluation of the model's ability to generalize to new data.

By following these rigorous training and validation protocols, we ensured that the models were well-prepared for accurate performance evaluation. The performance

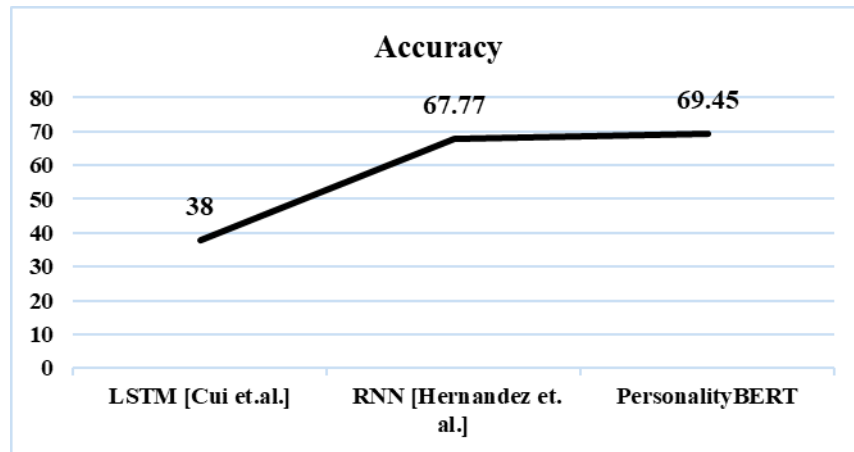
metrics presented in the following sections are based on models that have undergone this thorough training process.

## 7.2. Results of Research Objective 1

Research Objective 1 focused on developing and validating personality detection models that work across different languages.

### 7.2.1. Performance of PersonalityBERT

The results obtained by using the PersonalityBERT model are discussed in this section. It also mentions the hyper-parameter settings for the BERT model. State-of-the-art (SOTA) comparison is also done to validate the work. Performance of the proposed model is assessed using accuracy, recall, precision and f1 metrics. The dataset was divided in the ratio of 85:15 (where 85% was used for training the data and remaining 15% was used for testing purposes) using a ten-fold cross validation technique. This model was effectively used to classify the MBTI personality types from the English Kaggle\_MBTI dataset.



**Fig.7.1.** Accuracy Comparison with State-of-the-art

Fig 7.1 shows that PersonalityBERT achieves the highest accuracy among the models compared, with a score of 69.45%. This outperforms both the LSTM model, which has an accuracy of 38%, and the RNN model cited [Hernandez et al.], which has an accuracy of 67.77%. The application of PersonalityBERT has shown promising results in classifying complex personality traits using deep learning, highlighting the effectiveness of BERT-based models in processing nuanced language data.

### 7.2.2. Performance of HindiPersonalityNet

The evaluation metrics utilized to estimate the performance of the proposed HindiPersonalityNet model include accuracy, precision, recall, and F1-score. These metrics provide a comprehensive measure of the model's efficacy. While accuracy is

considered, F1-score is also incorporated to address the sensitivity of accuracy towards imbalanced datasets and to assess the classifier's performance in terms of specificity and sensitivity. To establish a benchmark, the results of the proposed model are compared with those of the state-of-the-art model. Table 7.1 showcases the performance results of the proposed model, providing a comprehensive overview of its effectiveness.

**Table 7.1.** Comparative Performance of Hindi Personality Detection Models

Model	Dataset	Accuracy	F1-score	Precision	Recall
<b>HindiPersonalityNet</b>	शख्सियत dataset	0.739	0.701	0.704	0.739
<b>KBSVE-P</b>	विशेष चरित्र_MBTI dataset	0.668	0.679	0.674	0.701

This model achieved the highest accuracy of 0.739 and recall of 0.739, indicating that it performed well in correctly identifying positive instances. It also had a relatively high F1-score of 0.701 and precision of 0.704. These results suggest that the HindiPersonalityNet model with GRU-BioWordVec embeddings is effective in capturing the patterns and features necessary for accurate predictions in the given task.

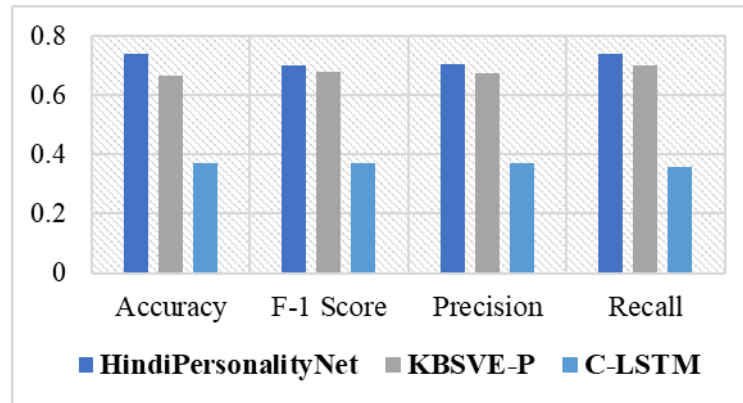
Table 7.2 presents a comparative performance analysis of two different models for Hindi and Bangla Personality Detection. The table provides information on the models, the datasets they were trained on, and the corresponding evaluation metrics including accuracy, F1-score, precision, and recall.

**Table 7.2.** Comparative Analysis of Hindi and Bangla Personality Detection

Model	Dataset	Accuracy	F1-score	Precision	Recall
<b>HindiPersonalityNet</b>	शख्सियत dataset	0.739	0.701	0.704	0.739
<b>C-LSTM [78]</b>	Bangla Personality Trait dataset	0.370	0.370	0.370	0.360

By comparing the performance metrics of the two models, it can be observed that the HindiPersonalityNet model achieved higher accuracy, F1-score, precision, and recall values compared to the C-LSTM model trained on the Bangla Personality Trait dataset. This analysis provides insights into the relative performance of the two models for personality detection in Hindi and Bangla languages.

Fig.7.2 presents the performance comparison of all 3 Indian language personality detection models.



**Fig.7.2.** Performance comparison of Indian Language personality detection

The results presented in the table 7.3 compare the performance of different deep learning models using BioWordVec.

**Table 7.3.** Comparative Analysis of various deep learning models with BioWordVec

Model	Accuracy	F1-score	Precision	Recall
<b>HindiPersonalityNet (GRU-BioWordVec)</b>	0.739	0.701	0.704	0.739
<b>CNN BioWordVec</b>	0.677	0.680	0.684	0.677
<b>BiLSTM BioWordVec</b>	0.679	0.675	0.674	0.679
<b>LSTM BioWordVec</b>	0.665	0.662	0.659	0.665

The CNN-based model achieved an accuracy of 0.677, which is lower than that of HindiPersonalityNet. However, it demonstrated a balanced F1-score of 0.680, precision of 0.684, and recall of 0.677. These results indicate that the CNN model performed reasonably well, but it fell short of the performance achieved by the HindiPersonalityNet model. Both BiLSTM BioWordVec and LSTM BioWordVec models exhibited similar performance. The BiLSTM BioWordVec achieved an accuracy of 0.679, an F1-score of 0.675, precision of 0.674, and recall of 0.679. On the other hand, the LSTM BioWordVec had an accuracy of 0.665, an F1-score of 0.662, precision of 0.659, and recall of 0.665. These models performed slightly lower than the HindiPersonalityNet and CNN models but still provided reasonably good results.

Table 7.4 presents a detailed comparative analysis of the HindiPersonalityNet model when combined with various deep learning models (CNN, LSTM, BiLSTM, GRU) and word embeddings (Crawl, fastText, GloVe, GoogleNews, PubMed, Wikipedia). The training and testing accuracies are reported for each combination, providing insights into how well the model performs on both the training and testing datasets.

**Table 7.4.** Accuracy comparison of embedding & deep learning combinations

Embedding + Deep Learning Model	Training Accuracy	Testing Accuracy
<b>Crawl CNN</b>	0.689	0.719
<b>Crawl LSTM</b>	0.735	0.724
<b>Crawl GRU</b>	0.742	0.721
<b>Crawl BiLSTM</b>	0.725	0.727
<b>fastText CNN</b>	0.757	0.724
<b>fastText LSTM</b>	0.821	0.722
<b>fastText GRU</b>	0.737	0.731
<b>fastText BiLSTM</b>	0.751	0.731
<b>Glove CNN</b>	0.707	0.727
<b>Glove LSTM</b>	0.749	0.725
<b>Glove GRU</b>	0.737	0.738
<b>Glove BiLSTM</b>	0.745	0.732
<b>GoogleNews CNN</b>	0.757	0.692
<b>GoogleNews LSTM</b>	0.853	0.698
<b>GoogleNews GRU</b>	0.766	0.725
<b>GoogleNews BiLSTM</b>	0.737	0.717
<b>PubMed CNN</b>	0.755	0.719
<b>PubMed LSTM</b>	0.784	0.727
<b>PubMed GRU</b>	0.752	0.731
<b>PubMed BiLSTM</b>	0.748	0.734
<b>Wikipedia CNN</b>	0.762	0.724
<b>Wikipedia LSTM</b>	0.752	0.712
<b>Wikipedia GRU</b>	0.764	0.721
<b>Wikipedia BiLSTM</b>	0.780	0.716

The table offers insights into the performance of different combinations of deep learning models and word embeddings for the HindiPersonalityNet model. It helps identify the most effective combinations in terms of accurately classifying instances and provides valuable information for selecting appropriate models and embeddings for similar tasks. By analysing the testing accuracies across different models, it is observed that fastText LSTM achieves the highest testing accuracy of 0.722, indicating its effectiveness in accurately classifying instances. The table also highlights the influence of different word embeddings on the model's performance. For instance, when combined with the GloVe embedding, the GRU model achieves a relatively high testing accuracy of 0.738. This suggests that the GloVe embedding provides valuable semantic information for the model to make accurate predictions.

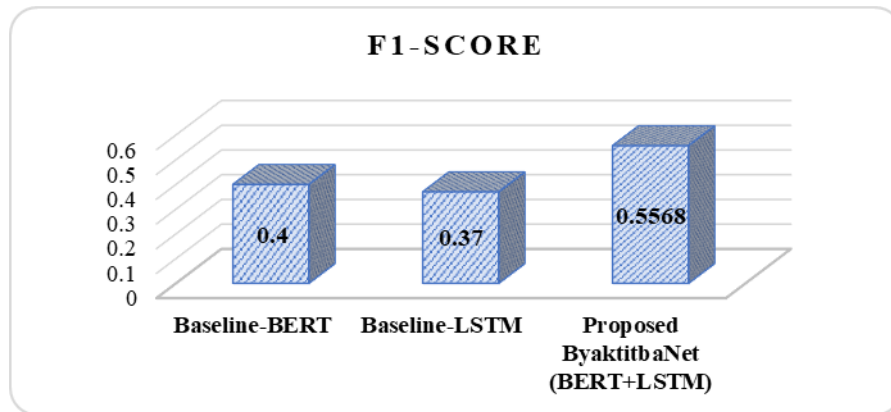
### 7.2.3. Performance of ByaktitbaNet

Accuracy, precision, recall and F1-score are used as the evaluation metrics for estimating the performance measure the efficacy of the proposed ByaktitbaNet model. F1-score is used along with accuracy because accuracy as a metric is quite sensitive towards the distribution of targeted values, and it is important to grasp the performance, specificity, and sensitivity of a classifier on an imbalanced dataset. Table 7.5 depicts the results of the proposed model along with baselines used.

**Table 7.5.** Performance of ByaktitbaNet Model

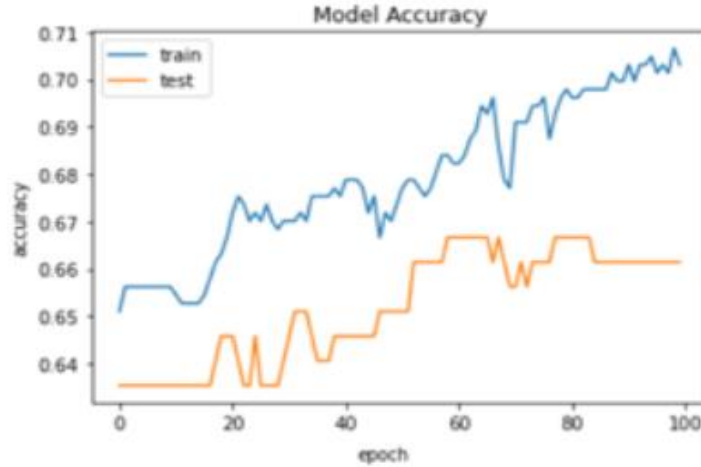
Model	Dataset	Accuracy	F1-score	Precision	Recall
<b>Baseline-BERT</b>	Training	0.5231	0.4190	0.4000	0.4120
<b>Test</b>	0.5454	0.4000	0.4000	0.4100	
<b>Baseline-LSTM</b>	Training	0.3600	0.3700	0.3600	0.3700
<b>Test</b>	0.3700	0.3700	0.3600	0.3700	
<b>Proposed ByaktitbaNet (BERT+LSTM)</b>	Training	0.6498	0.5119	0.4222	0.6498
<b>Test</b>	0.6849	0.5568	0.4691	0.6849	

Fig.7.3 clearly shows that the proposed ByaktitbaNet model outperforms the baseline BERT and LSTM models with an F1-score of 0.5568.



**Fig.7.3.** F1-score comparison

The fig.7.4 shows the exponential behaviour of the accuracy during the training and testing for the 100 epochs.



**Fig.7.4.** Model accuracy curve for 100 epochs

Table 7.6 shows the comparative results with a similar study.

**Table 7.6.** Performance comparison with another dataset

Model	Dataset	Labels	Accuracy	F1-score
<b>Convolutional-LSTM</b>	Bangla personality traits detection dataset	Big-5	0.37	0.37
<b>Proposed ByaktitbaNet (BERT+LSTM)</b>	Byaktitba (ব্যক্তিত্ব)	Introvert, Extrovert, Ambivert	0.6849	0.5568

Targeted for the Bangla dataset derived from "Baahubali 2", this model utilized BERT embeddings combined with LSTM to classify dialogues into introvert, extrovert, and ambivert categories. ByaktitbaNet achieves an accuracy of 0.6849, which suggests it is quite effective in classifying the personality traits within the dataset it was designed for. The accuracy indicates that the hybrid approach of ByaktitbaNet leveraging the contextual understanding capabilities of BERT and the sequence processing strengths of LSTM is well-suited for the task of personality detection in the Bangla language.

#### 7.2.4. Transformer-based Models for Personality\_Quotes

We have divided the results in two main tables to study the performance of deep learning models with static embeddings and the study the performance of transformer models on the English Personality\_Quotes dataset.

- **Performance of deep learning models**

Table 7.7 displays accuracies for different combinations of word embeddings and deep learning models. The accuracies are an indication of how well the models perform on a given task. Here's a breakdown of the tabular data:



- i. *Embedding Type:* The table lists different types of word embeddings used in the experiment: Common Crawl, FastText, GloVe, GoogleNews and PubMed.
- ii. *Deep Learning Models:* For each type of embedding, the experiment used different deep learning models: CNN (Convolutional Neural Network), LSTM (Long Short-Term Memory), GRU (Gated Recurrent Unit), and BiLSTM (Bidirectional LSTM).

**Table 7.7.** Accuracies for classical deep learning models with various static embeddings

<b>Embedding + Deep Learning Model</b>	<b>Accuracy</b>
<b>Common Crawl CNN</b>	0.728
<b>Common Crawl LSTM</b>	0.723
<b>Common Crawl GRU</b>	0.791
<b>Common Crawl BiLSTM</b>	0.752
<b>FastText CNN</b>	0.655
<b>FastText LSTM</b>	0.708
<b>FastText GRU</b>	0.786
<b>FastText BiLSTM</b>	0.708
<b>GloVe CNN</b>	0.762
<b>GloVe LSTM</b>	0.786
<b>GloVe GRU</b>	0.757
<b>GloVe BiLSTM</b>	0.776
<b>GoogleNews CNN</b>	0.752
<b>GoogleNews LSTM</b>	0.713
<b>GoogleNews GRU</b>	0.776
<b>GoogleNews BiLSTM</b>	0.762
<b>PubMed CNN</b>	0.737
<b>PubMed LSTM</b>	0.708
<b>PubMed GRU</b>	0.713
<b>PubMed BiLSTM</b>	0.723

Following observations can be made from the results shown in the table:

- a. *The Best Performing Combination:* The experiment indicates that the "Common Crawl" embeddings combined with the GRU model achieved the highest accuracy of 0.791. This suggests that for this particular dataset and task, using Common Crawl embeddings with the GRU model resulted in the most accurate predictions.
- b. *Embedding Comparison:* Among the different types of embeddings, the "GloVe" embeddings show better results across a variety of models (CNN,

LSTM, GRU, BiLSTM). This indicates that Glove embeddings are well-suited for this specific task compared to the other embeddings like fastText, GoogleNews and PubMed.

- c. **Model Performance:** Among the deep learning models, the results vary. GRU models generally perform well, but LSTM and BiLSTM models also show competitive performance. CNN models generally have slightly lower test accuracies compared to the RNN-based models (LSTM, GRU, BiLSTM).
- d. **Embedding Comparison with Models:** Combining the observations on embeddings and models, the best-performing combination is Common Crawl embeddings with the GRU model, achieving an accuracy of 0.791.

In conclusion, the results demonstrate that the choice of word embeddings and deep learning models significantly impacts the performance on the given task. Common Crawl embeddings with the GRU model stand out as the top-performing combination, while GloVe embeddings consistently perform well across multiple models.

- **Performance of transformer-based models**

The table 7.8 presents the performance metrics of transformer-based models for personality detection. The metrics include accuracy, F1 score, recall, precision, and the Matthews Correlation Coefficient (MCC).

**Table 7.8.** Performance of transformer-based models on *personality\_quotes* dataset

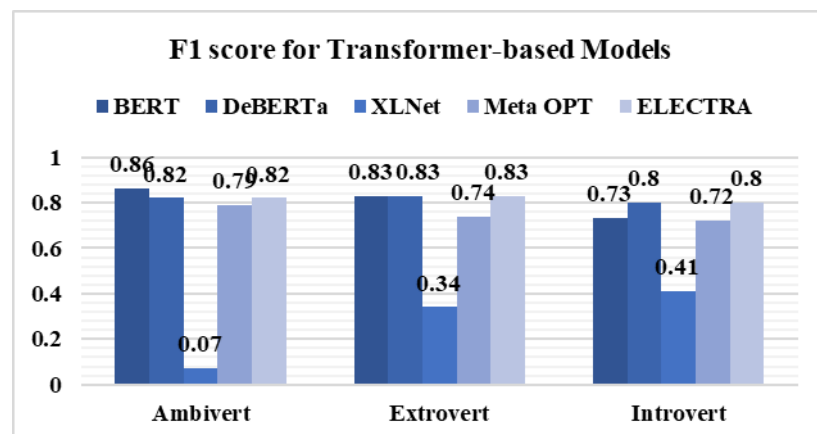
Model	Accuracy	F1_Score	Recall	Precision	MCC
<b>ELECTRA</b>	0.8106	0.8106	0.8106	0.8106	0.7195
<b>DeBERTa</b>	0.8106	0.8106	0.8106	0.8106	0.7156
<b>BERT</b>	0.8058	0.8058	0.8058	0.8058	0.7162
<b>Meta OPT</b>	0.7427	0.7427	0.7427	0.7427	0.6101
<b>XLNet</b>	0.3300	0.3300	0.3300	0.3300	-0.0578

The summarized analysis of the observations is as follows:

- **ELECTRA, DeBERTa, and BERT:** These models exhibit very similar and high levels of performance across all evaluated metrics, including accuracy, F1 score, recall, precision, and MCC. Their performance is consistent, indicating that they are well-suited for the task of personality detection.

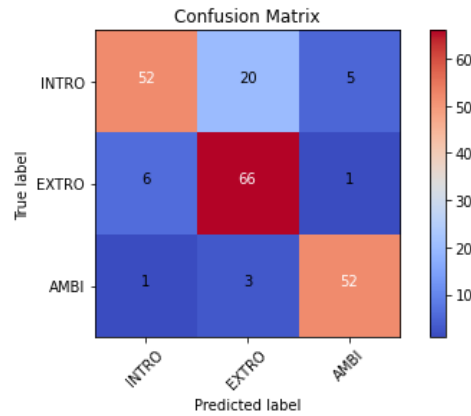
- **Meta OPT:** This model has slightly lower performance compared to ELECTRA, DeBERTa, and BERT. It has an accuracy of 0.7427 and an MCC of 0.6101. While its performance is still notable, it lags behind the others in terms of accuracy, F1 score, recall, precision, and MCC.
- **XLNet:** XLNet has the lowest performance among all models. Its accuracy, F1 score, recall, and precision are notably lower compared to the other models. The negative MCC value (-0.0578) suggests that its predictions are not consistent and might even be worse than random chance.

In summary, the observations highlight the performance differences between these transformer-based models on the personality detection task being evaluated. ELECTRA, DeBERTa, and BERT exhibit very similar and high levels of performance, Meta OPT has slightly lower performance, and XLNet performs the least effectively among the models, showing particularly poor performance according to the provided metrics. Figure 7.5 gives an overview of the F-1 scores achieved by the transformer-based models.



**Fig.7.5.** Comparison of f1 scores of transformer-based models

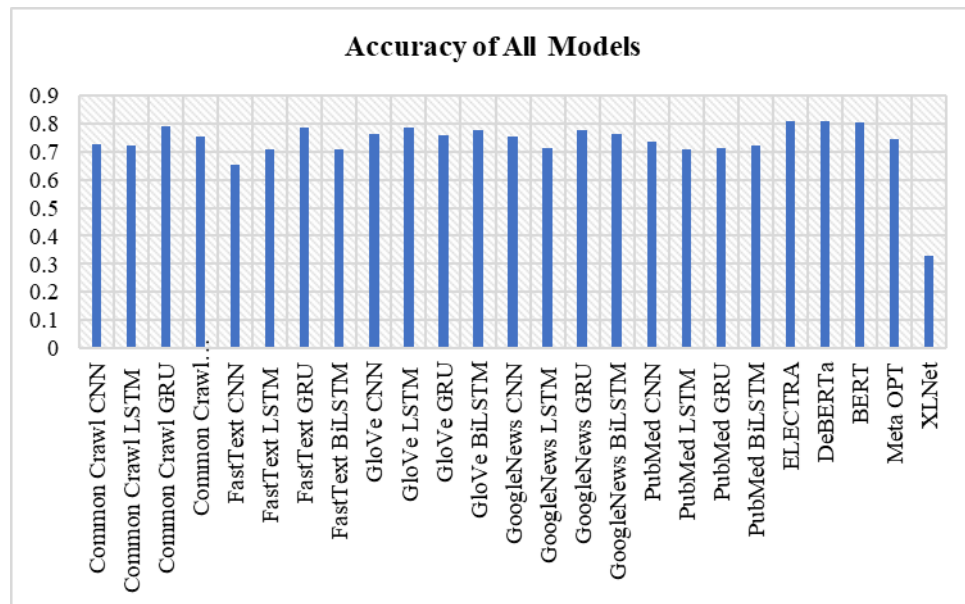
The confusion matrix for the top-performing model, ELECTRA is given in figure 7.6. The given confusion matrix represents a classification model's performance for three personality types: "Introvert," "Extrovert," and "Ambivert" on the test data. The training-test set split was taken as 80-20. Notably, it accurately predicted 52 instances for each of these classes. While there were various correct predictions, false positives occurred, including 7 Introverts misclassified as Introvert and 40 Extroverts mislabelled as Introvert. False negatives were also present, such as 25 instances of Introverts incorrectly identified as other classes. The matrix underscores the model's strengths and weaknesses in classifying the different personality types.



**Fig.7.6.** Confusion matrix for the top-performing model, ELECTRA

- **Comparison of the performance of classical deep learning models with transformer-based models**

Among the models, ELECTRA, DeBERTa, and BERT exhibit top performance with accuracies of 0.8106, showcasing the power of pretrained transformers. GRU and LSTM variants with various embeddings achieve competitive results, notably "Common Crawl GRU"(0.791), "GloVe LSTM"(0.786), and "GoogleNews GRU"(0.776). Notably, "XLNet" lags with an accuracy of 0.33. The diversity in results highlights the impact of embedding choice and architecture on outcomes, reinforcing the need for context-driven model selection. Figure 7.7 depicts the comparison of test accuracy of all models studied in this work.



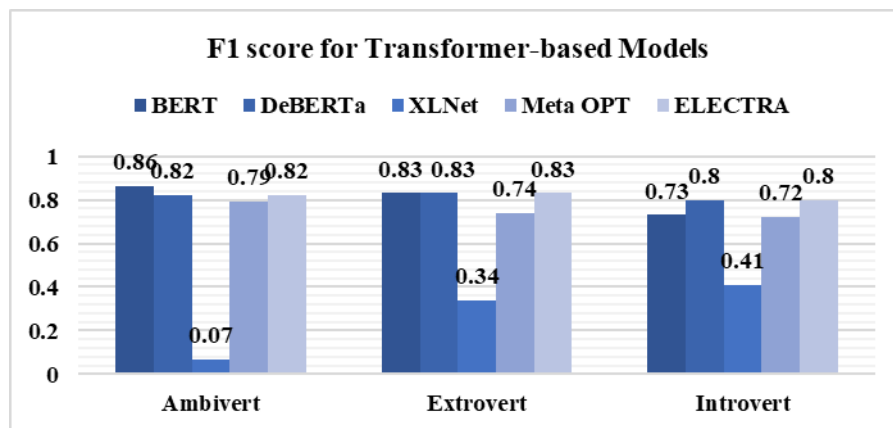
**Fig.7.7.** Comparison of accuracy of all models

Combining the observations on embeddings and models, the best-performing combination is Common Crawl embeddings with the GRU model, achieving an accuracy of 0.791. In conclusion, the results demonstrate that the choice of word embeddings and deep learning models significantly impacts the performance on the given task. Common Crawl embeddings with the GRU model stand out as the top-performing combination, while GloVe embeddings consistently perform well across multiple models. The table 7.9 presents the performance metrics of transformer-based models for personality detection.

**Table 7.9.** Performance of Transformer-Based Models on *Personality\_Quotes* Dataset

Model	Accuracy	F1_Score	Recall	Precision	MCC
<b>ELECTRA</b>	0.8106	0.8106	0.8106	0.8106	0.7195
<b>DeBERTa</b>	0.8106	0.8106	0.8106	0.8106	0.7156
<b>BERT</b>	0.8058	0.8058	0.8058	0.8058	0.7162
<b>Meta OPT</b>	0.7427	0.7427	0.7427	0.7427	0.6101
<b>XLNet</b>	0.3300	0.3300	0.3300	0.3300	-0.0578

Evidently, ELECTRA, DeBERTa, and BERT exhibit very similar and high levels of performance, Meta OPT has slightly lower performance, and XLNet performs the least effectively among the models, showing particularly poor performance according to the provided metrics. Figure 7.8 gives an overview of the F-1 scores achieved by the transformer-based models.



**Fig.7.8.** Comparison of F1 Scores of Transformer-Based Models

### 7.3. Results for Research Objective 2

Research Objective 2 enriched personality detection models by integrating deep psychological theories and psychometric methods, correlating constructs like parenting styles with MBTI types using data from low-resource linguistic datasets.

### 7.3.1. Model: KBSVE-P

We evaluate and observe the results of the KBSVE-P model on the South Asia low-resource Hindi Language dataset.

- **Performance of Kernel models: Kaggle\_MBTI**

All the personality traits have been studied over the four kernels. The performance of all the SVM kernels is evaluated using the six metrics are tabulated in Table 7.10. F1-Score is evaluated above accuracy because the accuracy metric is more delicate toward the distribution of target values; thus, it is important to capture a classifier's specificity, performance, and sensitivity on an imbalanced dataset.

**Table 7.10.** Performance of different SVM kernels: Kaggle\_MBTI

Kernels	Characteristic	Accuracy (%)	Precision (%)	Recall (%)	F1- Score (%)	ROC-AUC (%)	MCC( %)	Training Time (in sec)	Confusion Matrix
<b>RBF</b>	I/E	85.62	86.20	96.97	91.27	71.74	54.3	132.5	[227, 261], [51, 1630]
	N/S	88.24	88.81	98.82	93.55	60.32	36.2	111.3	[65, 233], [22, 1849]
	T/F	85.52	83.69	84.54	84.11	85.44	70.8	144.1	[1024, 162], [152, 831]
	J/P	79.48	80.90	63.10	70.90	76.66	56.4	168.1	[1182, 128], [317, 542]
<b>LINEAR</b>	<b>I/E</b>	86.17	87.71	95.54	91.46	74.72	56.9	82.9	[263, 225], [75, 1606]
	<b>N/S</b>	89.67	90.63	98.18	94.25	67.21	47.9	60.5	[108, 190], [34, 1837]
	<b>T/F</b>	84.	82.87	82.71	82.79	84.27	68.5	88.1	[1018, 68],

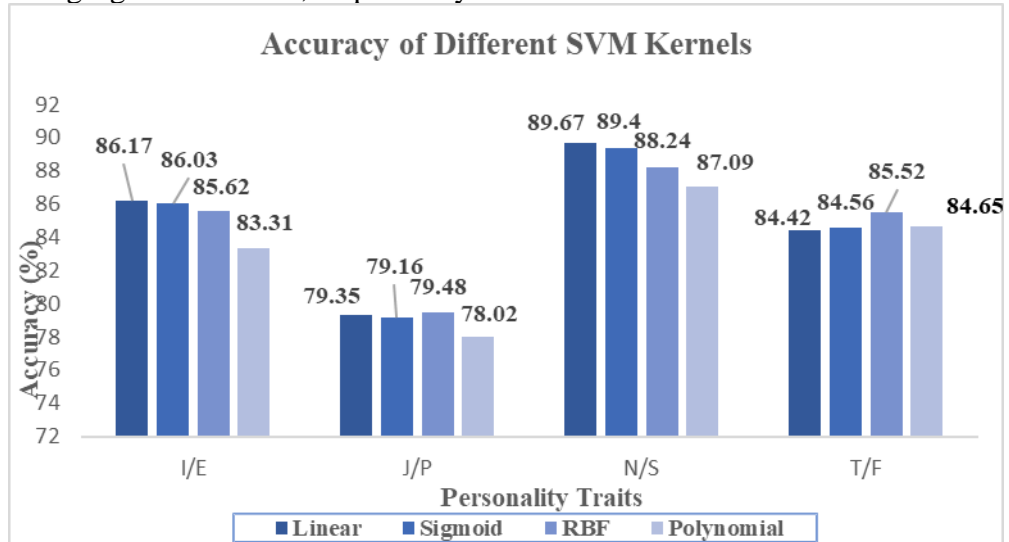
		42							[170, 813]
	<b>J/P</b>	79.35	78.42	66.01	71.68	77.05	56.1	109.1	[1154,156], [292,567]
<b>SIGMOID</b>	I/E	86.03	87.49	95.66	91.39	74.26	56.3	81.8	[258, 230], [73, 1608]
	N/S	89.40	90.20	98.40	94.12	65.64	45.7	58.6	[98, 200], [30, 1841]
	T/F	84.56	82.93	83.01	82.97	84.42	68.8	86.4	[1018,168], [167,816]
	J/P	79.16	78.38	65.42	71.32	76.80	55.7	107.5	[1155,155], [297,562]
<b>POLYNOMIAL</b>	I/E	83.31	83.36	98.04	90.10	65.31	45.1	201.4	[159, 329], [33, 1648]
	N/S	87.09	87.47	99.25	92.99	54.99	23.8	190.8	[32, 266], [14, 1857]
	T/F	84.65	81.99	84.74	83.34	84.66	69.1	207.1	[1003,183], [150,833]
	J/P	78.61	85.08	55.76	67.37	74.68	55.1	209.3	[1226, 84], [380, 479]

I/E: Introvert-Extrovert; N/S: iNtution-Sensing; T/F: Thinking-Feeling; and J/P: Judgemental-Perceiving.

From Table 7.10, the following observations are made. First, the accuracy of the linear kernel is found to be highest in the case of I/E and N/S personality traits. Whereas in the case of T/F and J/P personality traits, the accuracy of RBF is highest. Second, the F1 score of the linear kernel turned out to be the highest in the case of I/E, N/S, and J/P personality traits. On the other hand, RBF shows the highest F1 score in the case of the T/F personality trait. Third, the ROC-AUC of the linear kernel is found to be highest in the case of I/E, N/S, and J/P personality traits. Whereas in the case of T/F personality trait, the ROC -AUC of RBF is highest. Fourth, the MCC of the linear kernel turned out to be the highest in the case of I/E

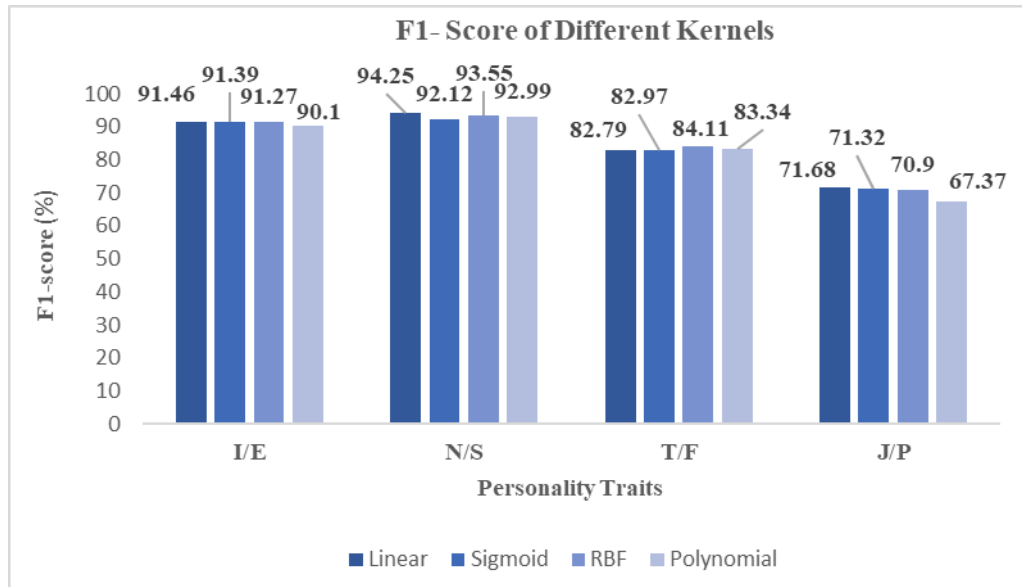
and N/S personality traits. On the other hand, RBF shows the highest MCC in the case of T/F and J/P personality traits. At last, from the above four observations, it can be said that the linear kernel performs better in most of the performance matrices when compared to other kernels.

The corresponding visualization of accuracies and F1-scores are shown in the following fig. 7.9 and 7.10, respectively.



**Fig. 7.9.** Accuracies of different SVM kernels: Kaggle\_MBTI

Fig.7.9 shows that the accuracy of the linear kernel is highest in the case of I/E (86.17) and N/S (89.67) personality traits. On the other hand, the accuracy of RBF is highest in the case of J/P (79.48) and T/F (85.52).



**Fig. 7.10.** F1-Score of different SVM kernels: Kaggle\_MBTI



Further from fig. 7.10, we can infer that the F1 score of the linear kernel is highest in the case of I/E (91.46), N/S (94.25), and J/P (71.68) personality traits, whereas in the case of T/F personality trait, RBF shows the highest F1 score. Hence, we conclude that the linear kernel performs better in 3 out of 4 cases as compared to the other three kernels.

The results show that the linear kernel performs best among all kernels with average accuracy and F1 score of 84.7 and 84.9, respectively for the personality traits.

- **Performance of voting techniques: Kaggle\_MBTI**

The results obtained after training and testing SVM kernels over the MBTI personality traits aggregated with the three different voting techniques such as soft voting, hard voting, and weighted hard voting are depicted in Table 7.11.

**Table 7.11.** Performance of ensemble of SVM kernels aggregated with different voting techniques

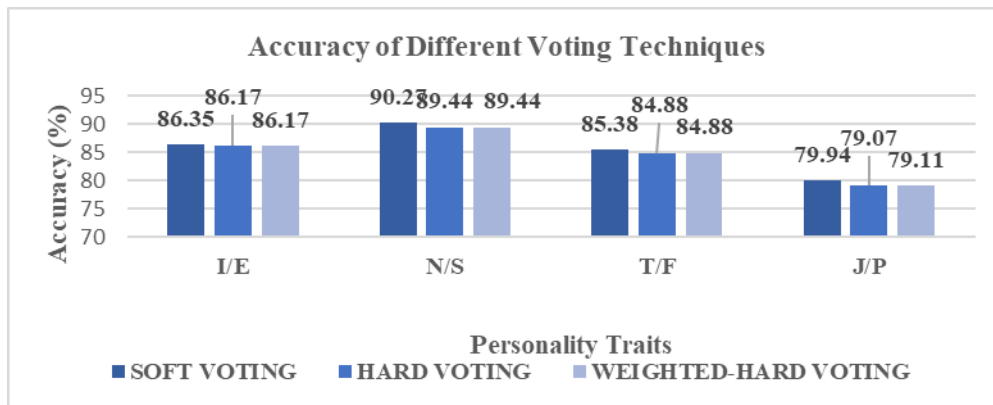
Voting	Characteristic	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)	ROC-AUC (%)	MCC (%)	Training Time (In Second)	Confusion Matrix
SOFT VOTING	IE	86.35	88.24	95.06	91.52	75.71	57.8	2264.6	[275, 213], [ 83, 1598]
	NS	90.27	91.58	97.70	94.54	70.66	52.5	1884.4	[130, 168], [43, 1828]
	TF	85.38	83.70	84.13	83.92	85.28	70.5	2380.1	[1025,161], [156,827]
	JP	79.94	78.80	67.52	72.73	77.81	57.4	2695.3	[1154,156], [279,580]

<b>HARD VOTING</b>	IE	86.17	87.59	95.72	91.47	74.50	56.8	507.1	[260, 228], [72, 1609]
	NS	89.44	90.25	98.40	94.14	65.81	46.1	408.8	[99, 199], [ 30, 1841]
	TF	84.88	84.01	82.30	83.14	84.66	69.4	512.7	[1032,154], [174,809]
	JP	79.07	80.82	61.82	70.05	76.10	55.5	597.3	[1184,126], [328, 531]
<b>WEIGH TED HARD VOTING</b>	IE	86.17	87.59	95.72	91.47	74.50	56.8	476.2	[260, 228], [ 72, 1609]
	NS	89.44	90.25	98.40	94.14	65.81	46.1	395.4	[99, 199], [ 30, 1841]
	TF	84.88	84.01	82.30	83.14	84.66	69.4	501.8	[1032,154], [174, 809]
	JP	79.11	80.85	61.93	70.14	76.16	55.6	561.8	[1184,126], [327, 532]

I/E: Introvert-Extrovert; N/S: iNtution-Sensing; T/F: Thinking-Feeling and J/P: Judgemental-Perceiving.

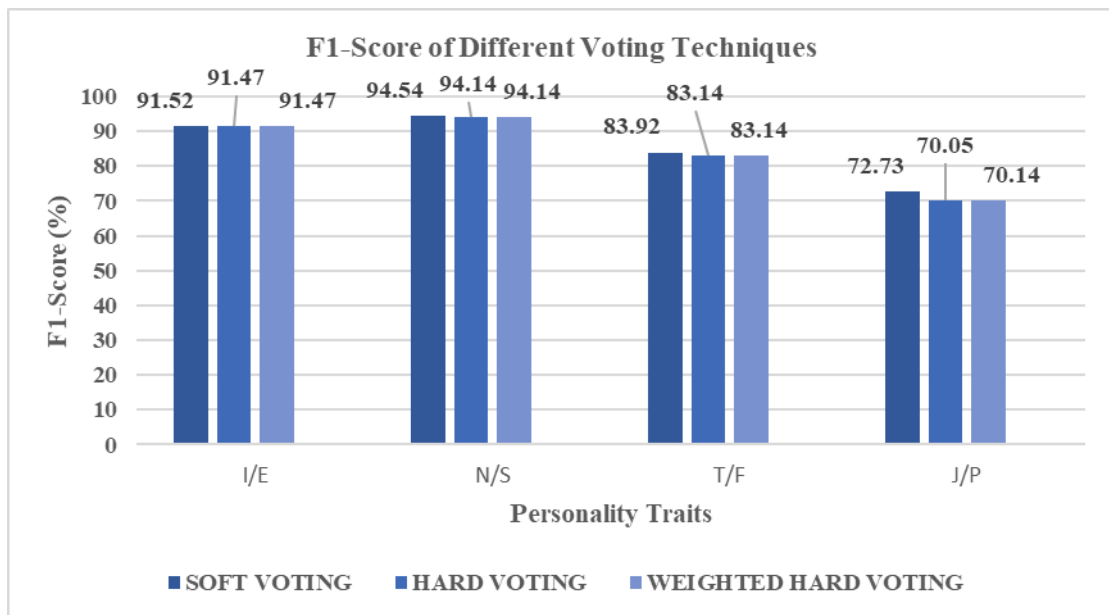
From Table 7.11, it can be inferred that soft voting outperforms the other voting techniques, such as hard voting and weighted hard voting, while considering most of the performance matrices, such as accuracy, F1 score, ROC-AUC, and MCC. Here, another vital point to note is that hard voting and weighted hard voting perform almost identical across all performance matrices.

The corresponding visualizations of accuracy and F1-score are shown below in fig. 7.11 and fig. 7.12, respectively.



**Fig. 7.11.** Accuracies of voting techniques: Kaggle\_MBTI

From fig. 7.11, we can conclude that the average accuracy of the soft voting classifier is 85.48, which is higher than the other voting classifiers, such as the hard voting and the weighted hard voting classifier. Here average accuracy of hard voting and weighted hard voting classifier is 84.89 and 84.9, respectively. Moreover, another critical point to note is that the accuracy of the soft voting classifier is better than other voting classifiers in the case of all personality traits.



**Fig. 7.12.** F1-Score of different Voting techniques: Kaggle\_MBTI

Further from fig. 7.12, we can infer that the F1 score of the soft voting classifier is the highest for all personality traits when compared to that for hard voting and weighted hard voting classifiers. Here average F1 score for the soft voting classifier

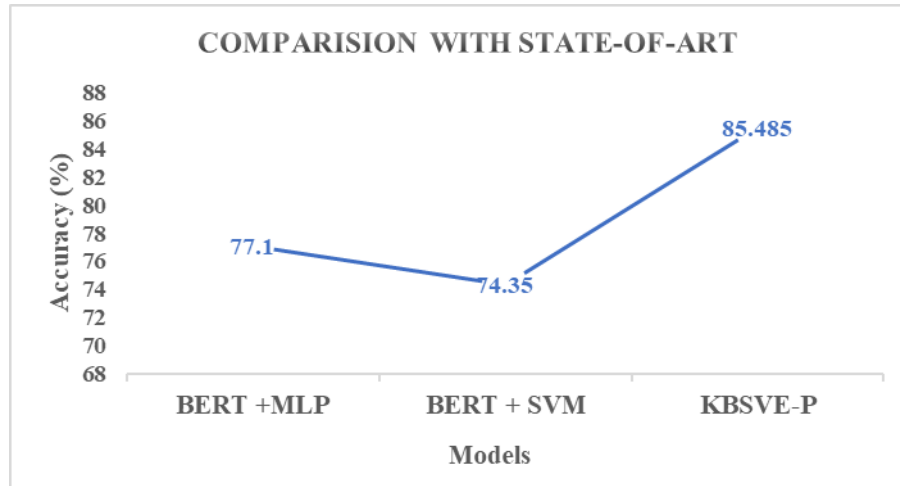
is 85.67, while it is 84.7 and 84.72 for hard voting and weighted hard voting classifiers.

- **Comparison with existing works: Kaggle\_MBTI**

The results provided by the proposed kernel-based ensemble model, KBSVE-P, are compared with the state-of-the-art on Kaggle-MBTI. The proposed model outperformed the existing works. Table 7.12 results show an improved performance of 9.4%, 4.6%, 12.19%, and 18.95% for I/E, N/S, T/F, and J/P MBTI personality traits, respectively.

**Table 7.12.** Comparison with the State-of-Art (Accuracy): Kaggle\_MBTI

Model used	I/E	N/S	T/F	J/P
<b>BERT + MLP</b>	78.8	86.3	76.1	67.2
<b>BERT + SVM</b>	77.0	86.2	73.7	60.5
<b>Proposed KBSVE-P</b>	<b>86.35</b>	<b>90.27</b>	<b>85.38</b>	<b>79.94</b>



**Fig. 7.13.** Comparison with the existing work in Kaggle\_MBTI

Fig. 7.13 shows the accuracy comparison of models, and the proposed KBSVE-P model outperforms the existing models with an average accuracy of 85.48 for the personality traits.

- **Evaluation of KBSVE-P on विशेषचरित्र\_\_MBTI dataset**

The proposed KBSVE-P model is also evaluated on the South Asian, Indian low-resource Hindi language dataset, विशेषचरित्र\_\_MBTI dataset (vishesh charitr\_MBTI). Table 7.13 shows the accuracy of the model.

**Table 7.13.** KBSVE-P on विशेषचरित्र\_\_MBTI dataset

Dataset	Accuracy (%)	Precision (%)	Recall (%)	F1-score (%)	ROC_AUC (%)	MCC (%)
विशेष चरित्र__MBTI (vishesh charitr)	66.89	67.40	70.17	67.96	64	30.09

In summary, the KBSVE-P models performance varies with the choice of SVM kernels and voting techniques for different personality trait predictions. The RBF kernel generally provides high accuracy for the N/S and T/F dimensions, while Linear and Polynomial kernels are more effective for the I/E dimension. When it comes to voting classifiers, Soft Voting tends to yield higher accuracy across most personality traits, highlighting its effectiveness as an ensemble method in personality detection tasks. This suggests that Soft Voting can be a preferable approach when utilizing the KBSVE-P model for classifying MBTI personality types.

### 7.3.2. Performance of DM\_MBTI

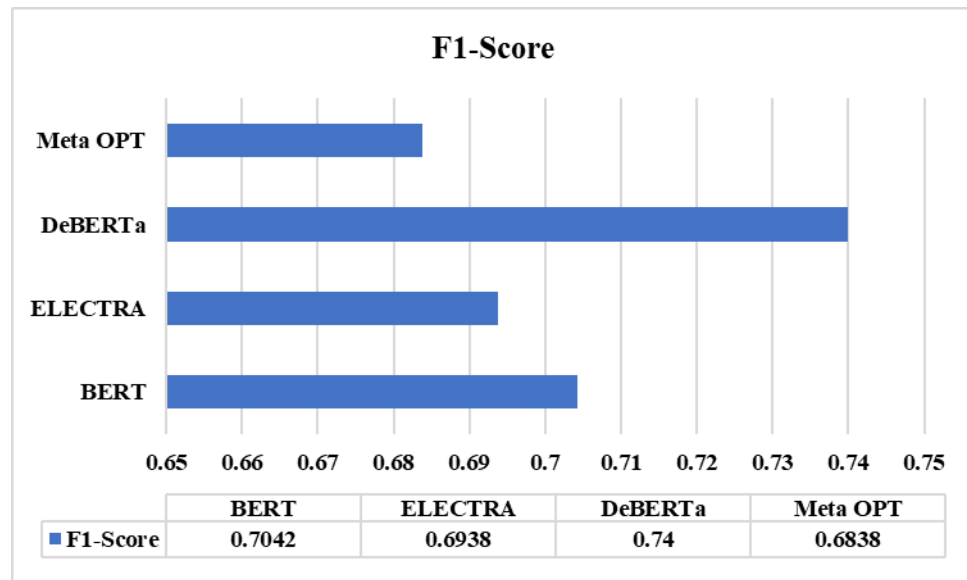
#### • Performance of Transformer-based Models

In Table 7.14, we present an array of performance metrics for transformer-based models used in MBTI decision-maker type recognition. These metrics encompass accuracy, F1 score, recall, precision, and the Matthews Correlation Coefficient (MCC). Accuracy quantifies the proportion of correctly predicted instances relative to the total instances, offering an overarching measure of model correctness. The F1 Score, a balanced metric, harmonizes precision (the capability to avoid mislabelling negative samples as positive) and recall (the ability to detect all positive samples). Recall assesses the model's proficiency in identifying positive cases by measuring the proportion of true positive instances among all actual positive instances. Precision, on the other hand, gauges the accuracy of predicted positive cases by calculating the proportion of true positive instances relative to all predicted positives. Lastly, MCC, the Matthews Correlation Coefficient, serves as a quality gauge for binary classifications, factoring in true and false positives and negatives. A higher MCC signifies superior model performance, with 1 signifying perfection, 0 indicating randomness, and -1 indicating complete discord between predictions and actual labels.

**Table 7.14.** Performance of all Transformer-based Models on DM-MBTI dataset

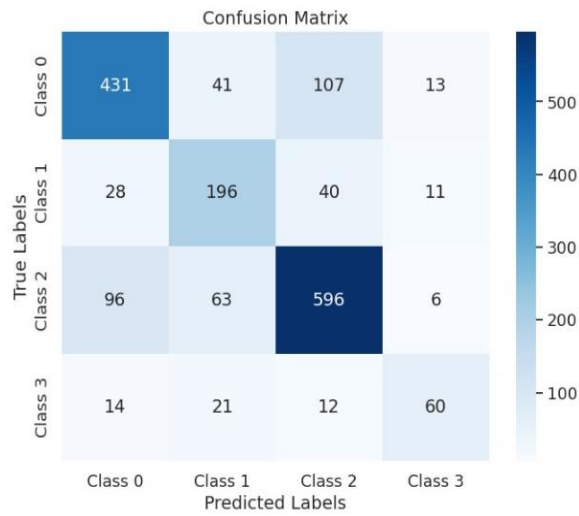
Model	Accuracy	Precision	Recall	F1-Score	MCC
BERT	0.7135	0.7226	0.7135	0.7042	0.5583
ELECTRA	0.6997	0.7002	0.6997	0.6938	0.5373
DeBERTa	0.7395	0.7426	0.7395	0.74	0.6085
Meta OPT	0.6893	0.6887	0.6893	0.6838	0.5207

BERT achieves a good level of accuracy, indicating its ability to correctly classify decision-maker types. Precision is also relatively high, which means that when BERT predicts a decision-maker type, it's often correct. However, the F1-Score and MCC values are slightly lower, suggesting a challenge in balancing precision and recall. This indicates that while BERT makes accurate predictions, it might not capture all decision-maker types, leading to some missed classifications. ELECTRA's performance is similar to BERT, with comparable accuracy, precision, and recall. The F1-Score and MCC values are also in the same range, suggesting a balanced but moderate performance. It makes accurate predictions but may not capture all decision-maker types. Meta OPT performs slightly lower across all metrics compared to the other models. While it maintains decent accuracy and precision, its F1-Score and MCC values are lower. This suggests that Meta OPT may struggle with capturing a comprehensive range of decision-maker types, potentially leading to some misclassifications. Fig. 7.14 visually depicts the comparison of F1 scores across all the evaluated Transformer-based models.



**Fig.7.14.** F1 comparison of Transformer-based Models

DeBERTa's superior performance in terms of accuracy, precision, recall, F1-Score, and MCC score, coupled with its balanced trade-off between precision and recall, makes it the best choice for Decision Maker MBTI recognition. It consistently provides accurate results while effectively capturing the diverse range of decision-maker types, making it a robust and reliable model for this specific task. Fig. 7.15 shows the confusion matrix for DeBERTa where Consensual is class 0, Consultative is class1, Deliberate is class 2 and Quick is class 3.



**Fig.7.15.** Confusion matrix for the best performing Transformer model, DeBERTa

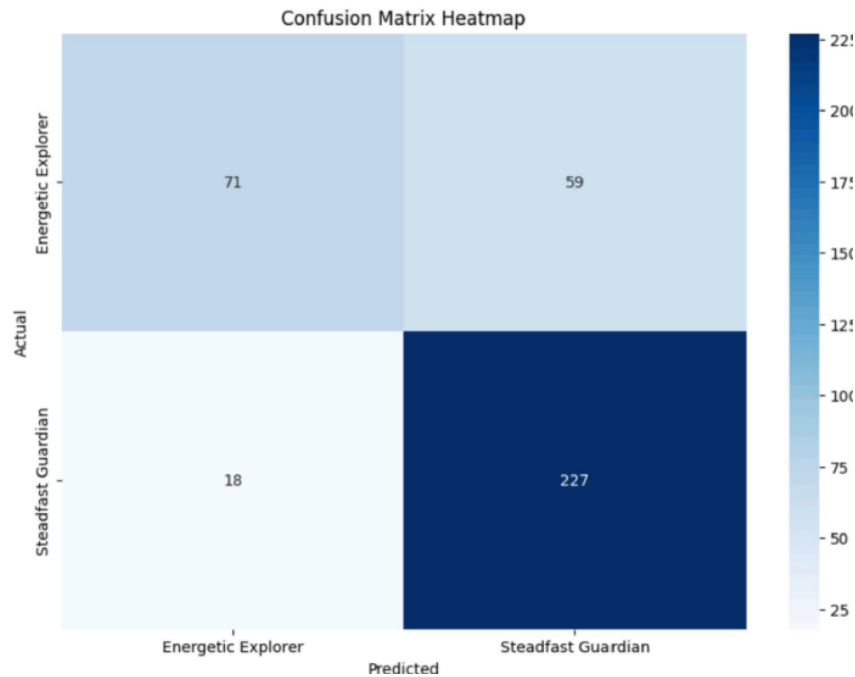
### 7.3.3. Performance of ParvarishNET

Table 7.15 displays the results of the ParvarishNET model which is the XGBoost Model

**Table 7.15.:** Performance of ParvarishNET

Model	Accuracy	F1-score	Recall	Precision
ParvarishNET (XGBoost)	<b>0.79</b>	<b>0.75</b>	<b>0.74</b>	<b>0.79</b>
mBERT	0.75	0.71	0.70	0.73





**Fig.7.16.:** Confusion matrix of the ParvarishNET

The above fig. 7.16 shows the confusion matrix of the best performing model, ParvarishNET, an XGBoost Model.

**Table 7.16:** Comparison of the results of parenting styles

Model	Parenting Style	Precision	Recall	F1-Score
mBERT	Energetic Explorers	0.69	0.53	0.6
	Steadfast Guardians	0.77	0.87	0.82
<b>ParvarishNET (XGBoost)</b>	<b>Energetic Explorers</b>	<b>0.8</b>	<b>0.55</b>	<b>0.65</b>
	<b>Steadfast Guardians</b>	<b>0.79</b>	<b>0.93</b>	<b>0.85</b>

The table 7.16 compares the performance of XGBoost, a tree-based ensemble model with mBERT, a transformer model. XGBM outperforms mBERT across both parenting styles, particularly in terms of precision, recall, and F1-Score. XGBoost's tree-based architecture excels at capturing complex, non-linear feature interactions, which helps reduce false positives (higher precision) and identify more true positives (higher recall), leading to a balanced F1-Score. mBERT, while strong in understanding contextual nuances due to its Transformer-based architecture, struggles with precision and recall, likely due to overfitting on less relevant features or failing to capture specific patterns critical to classification. XGBoost ensemble approach also enhances robustness and reduces variance, making it more effective for tasks requiring precise and balanced predictions across diverse categories.

Table 7.17 has the detailed ablation study of different models used over Parvarish dataset.

**Table 7.17:** Performance of different models over Parvarish Dataset

Model	Accuracy
BERT	66.13
DeBERTa	65.33
LSTM	66
RoBERTa	74
mBERT	75
<b>XGBoost</b>	<b>79</b>

Table 7.17 describes that XGBoost, with an accuracy of 79%, outperforms the other models, likely due to its strength in capturing complex, non-linear patterns through its ensemble tree-based approach. mBERT, achieving 75% accuracy, benefits from its extensive multilingual training, while RoBERTa's 74% accuracy reflects the effectiveness of its improved pre-training techniques. In contrast, BERT and LSTM show similar performance with accuracies of 66.13% and 66%, respectively, indicating BERT's strong contextual understanding and LSTM's capability in handling sequential data. DeBERTa slightly underperforms with 65.33% accuracy, suggesting that its architectural enhancements did not significantly boost performance in Parvarish Dataset.

#### 7.4. Results for Research Objective 3

Research Objective 3 tackled profiling emotional dispositions like optimism and pessimism, assessing model applicability in practical domain.

##### 7.4.1. Performance of Transformer-based for Manobhav Dataset

This section encompasses a detailed analysis of model performance, their potential, and the hyperparameters crucial for their effectiveness. The evaluation provides valuable insights into model choice and improvement strategies, fostering a deeper understanding of the models' capabilities in text-based attitude prediction tasks. A concise summary of the limitations of the study is also presented.

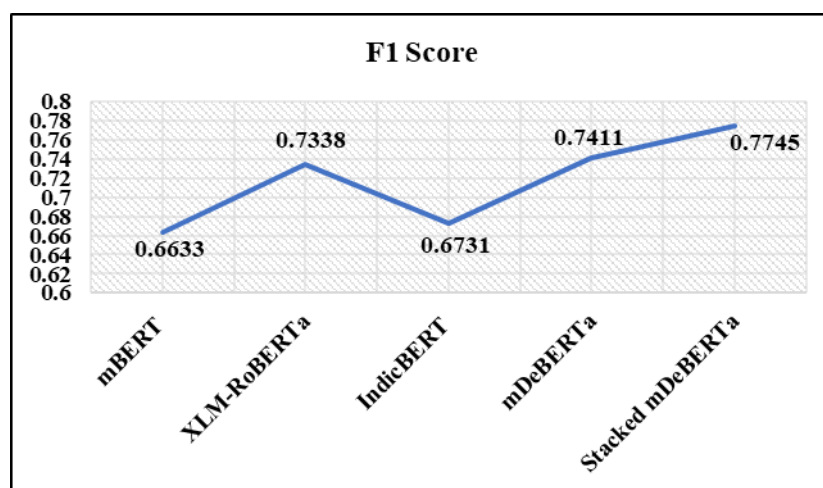
- **Performance of Transformer-Based Models on मनोभाव (pronounced as Manobhav) dataset**

We provide a comprehensive analysis of the performance of various Transformer-based models on the 'मनोभाव' dataset. A comparative evaluation based on metrics including Accuracy, Precision, Recall, F1 Score, and ROC AUC Score is presented in Table 7.18. The results highlight the effectiveness of each model in capturing the nuances of the dataset.

**Table 7.18.:** Performance of all Transformer-based Models on मनोभाव dataset

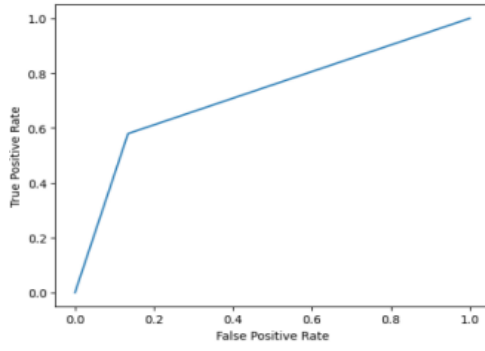
Models	Accuracy	Precision	Recall	F1 Score	ROC AUC Score
mBERT	0.6712	0.6580	0.6712	0.6633	0.5911
XLM-RoBERTa	0.7474	0.7332	0.7474	0.7338	0.6618
IndicBERT	0.6851	0.6672	0.6851	0.6731	0.5979
mDeBERTa	0.7543	0.7411	0.7543	0.7411	0.6700
Stacked mDeBERTa	0.7785	0.7726	0.7785	0.7745	0.7226

Furthermore, Fig. 7.17 visually depicts the comparison of F1 scores across all the evaluated Transformer-based models.

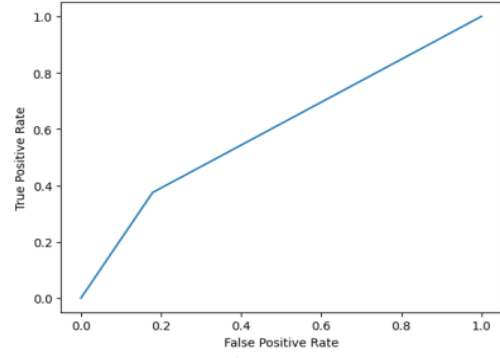


**Fig.7.17.:** F1 comparison of Transformer-based Models

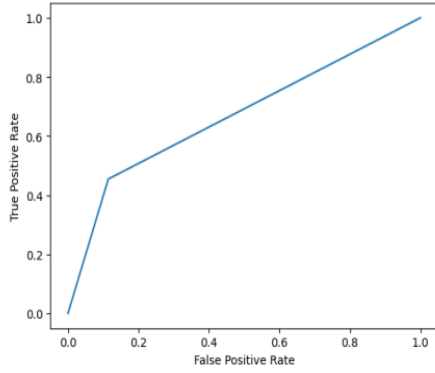
It is evident that the proposed Stacked mDeBERTa model outperforms the other models in terms of accuracy, precision, recall, F1 score, and ROC AUC score. It demonstrates a well-rounded performance in identifying positive cases while maintaining a good balance between avoiding false positives and capturing true positives. The high ROC AUC score (0.7226) indicates that the model is effective at ranking positive cases higher than negative cases across a range of thresholds. This implies a good ability to distinguish between the two classes. The ROC curves of all five models are shown in fig. 7.18 to fig.7.22.



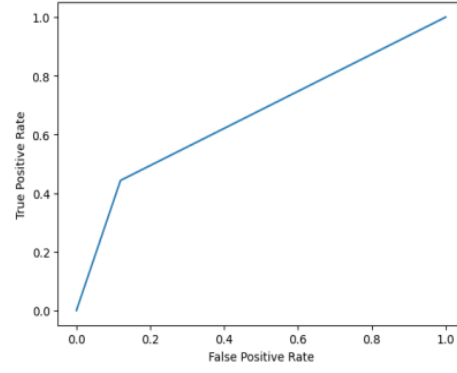
**Fig.7.18.** ROC curve for Stacked mDeBERTa  
IndicBERT



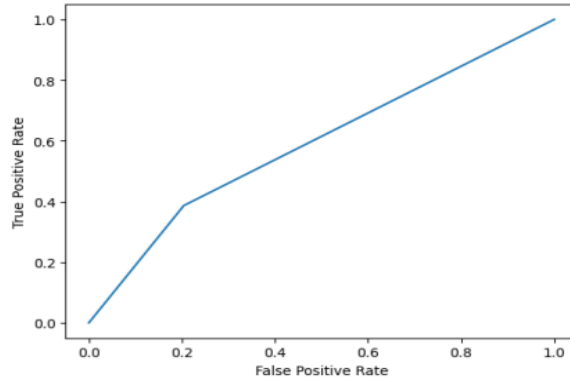
**Fig.7.19.** ROC curve for



**Fig.7.20.** ROC curve for mDeBERTa

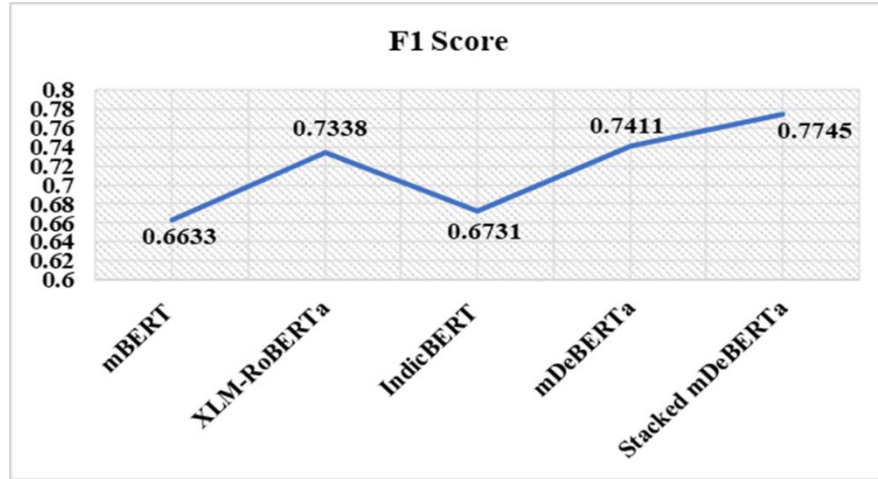


**Fig.7.21.** ROC curve for XLM-RoBERTa



**Fig.7.22.** ROC curve for mBERT

Focused on profiling emotional attitudes like optimism and pessimism in Hindi, this model achieved high accuracy using stacked mDeBERTa configurations. Figure 7.23 showcases the F1 scores of various transformers used on the Manobhav dataset.



**Fig. 7.23.** F1 score comparison of transformer-based models

Based on the graph for the Manobhav dataset mBERT model scored an F1 of approximately 0.6633, indicating a moderate level of precision and recall balance. XLM-RoBERTa model shows an improvement with an F1 score of around 0.7338. IndicBERT has a slight dip in performance compared to XLM-RoBERTa, with an F1 score close to 0.6731. mDeBERTa model has a significant improvement, marking an F1 score of about 0.7411. The most impressive results come from the Stacked mDeBERTa model, achieving the highest F1 score of approximately 0.7745 on the graph.

From these results, it's clear that the stacked mDeBERTa model outperforms the others on the Manobhav dataset. The high F1 score suggests that it has a strong balance between precision and recall, making it reliable for detecting emotional attitudes in the Hindi language. The performance increase from mBERT to stacked mDeBERTa illustrates the benefits of stacking models and potentially using ensemble techniques to enhance predictive capabilities.

#### 7.4.2. Performance of Multimodal fusion architecture

A comprehensive analysis of the experimental results obtained from the multimodal models used for predicting personality traits is discussed. Evaluation metrics such as accuracy, F1-score, precision, and recall are employed to assess the performance of these models. The discussion highlights the model performance across different data modalities and provides insights into the effectiveness of each approach. The following hyperparameters (table 7.19) have been tuned while training the models:

**Table 7.19.** Model Hyperparameters

Model	Hyperparameter
BERT	Learning Rate = 5e-5
DeBERTa	Learning Rate = 1e-4
InceptionV3	Learning rate = 0.03 Alternate ReLU and Tanh activation functions for Dense layers SGD optimizer with 0.9 momentum instead of Adam allows the gradient to not get stuck at a local minimum easily
EfficientResNet	Learning rate = 0.01 initially, with an exponential decay with weight decay = 0.3 over 500 steps (approximately 2.5 epochs) Sigmoid dense layer followed by a linear regression layer for output. Weighted fusion of convolutional outputs from frames, audio data and transformers output of transcript with weights 0.1 for each frame, 0.4 for audio and 0.3 for transcript
XceptionResNet	Learning rate = 0.001 Weighted fusion of convolutional outputs from frames, audio data and transformers output of transcript with weights 0.2 for each frame, 0.3 for audio and 0.1 for transcript

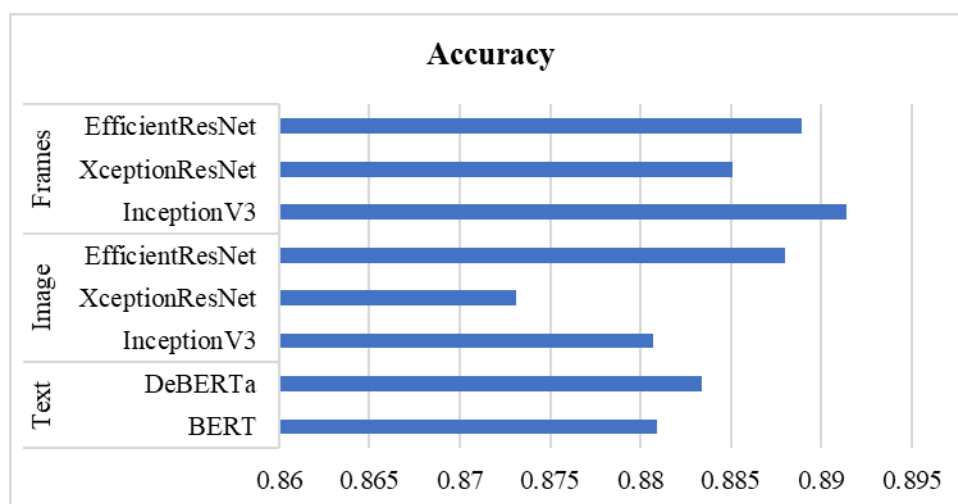
- **Performance across modalities**

We begin by analysing the performance of our models across various data modalities. The table 7.20 below summarizes the results for different models and modalities:

**Table 7.20.** Performance across modalities

Modalities	Model	Accuracy	R2 Score	Mean Squared Error	Root Mean Squared Error
<b>Text</b>	BERT	0.8809	-2.60	0.0216	0.1469
	DeBERTa	0.8834	-8.31 x 10 <sup>9</sup>	0.0215	0.1467
<b>Image</b>	InceptionV3	0.8807	-1.62%	0.0219	0.1478
	EfficientResNet	0.8880	9.97%	0.0193	0.1390
	XceptionResNet	0.8731	-14.46%	0.0247	0.1569
<b>Frames</b>	InceptionV3	0.8914	15.37%	0.0182	0.1347
	EfficientResNet	0.8889	12.12%	0.0189	0.1374
	XceptionResNet	0.8851	4.64%	0.0205	0.1431
<b>Fusion/ Decision models</b>	InceptionV3 + BERT	0.8905	13.05%	0.0187	0.1366
	EfficientResNet + BERT	0.8876	10.55%	0.0194	0.1394
	XceptionResNet + BERT	0.9212	54.49%	0.0098	0.0992

From table 7.20 we observe that the performance of the models varies across modalities. For text, DeBERTa performs poorly in terms of R2 score, while for audio and video, EfficientResNet and InceptionV3 perform well, respectively. When considering all modalities, XceptionResNet + BERT stands out with the highest accuracy, a high R2 score, and the lowest mean squared error and root mean squared error. It appears to be the best-performing model for the task involving all modalities. The bar chart in fig.7.24 illustrates the accuracy of various deep learning models in processing text, image, and frame data modalities for personality analysis, with EfficientResNet showing a leading performance.



**Fig.7.24.** Model Comparison for Personality Analysis: Accuracy of Multimodal Data Processing

#### • Results on Text Modality

The results for two different models namely BERT and DeBERTa in terms of average accuracy and accuracy on various personality traits are in table 7.21 given below.

**Table 7.21.** Performance of Text Modality Models

Model	Avg. Accuracy	Accuracy on Various Traits				
		Extraversion	Openness	Agreeableness	Neuroticism	Conscientiousness
<b>BERT</b>	88.09%	89.24%	87.25%	87.85%	87.81%	88.28%
<b>DeBERTa</b>	88.34%	89.49%	88.23%	87.77%	87.70%	88.52%

Both models demonstrate high average accuracy, with BERT scoring 88.09% and DeBERTa scoring 88.34%. This suggests that both models perform well on the overall task of personality trait prediction from text data. Analysing the trait-wise



accuracies, we find that Extraversion is the easiest trait to identify using the transcripts of the videos. This implies that individuals' spoken words can effectively predict whether they are extroverts or introverts. However, the other traits are relatively more challenging to identify from text data. Notably, Openness proves to be the most difficult to predict, likely requiring the assessment of body language and non-verbal cues.

- **Results on Image Modality Generated from Audio Files**

Moving on to the image modality derived from audio files, we summarize the results in Table 7.22:

**Table 7.22.** Performance of Audio-generated Images Models

Model	Avg. Accuracy	Accuracy on Various Traits				
		Extraversion	Openness	Agreeableness	Neuroticism	Conscientiousness
<b>InceptionV3</b>	88.07%	89.42%	88.47%	87.40%	87.67%	87.38%
<b>EfficientResNet</b>	88.80%	89.98%	89.23%	88.29%	88.37%	88.13%
<b>XceptionResNet</b>	87.31%	88.70%	87.76%	86.36%	87.01%	86.73%

In this modality, Extraversion remains the most predictable trait, closely followed by Openness. Conscientiousness appears to be the most challenging trait to identify. EfficientResNet appears to perform the best overall, with the highest average accuracy and the highest accuracy on most traits.

- **Results on Frames Generated from Video Files**

For frames generated from video files, the results are summarized in Table 7.23:

**Table 7.23.** Performance of Video frames

Model	Avg. Accuracy	Accuracy on Various Traits				
		Extraversion	Openness	Agreeableness	Neuroticism	Conscientiousness
<b>InceptionV3</b>	89.14%	89.93%	89.66%	88.63%	88.89%	88.59%
<b>EfficientResNet</b>	88.89%	89.99%	89.28%	88.35%	88.61%	88.21%
<b>XceptionResNet</b>	88.51%	89.54%	88.71%	87.78%	88.64%	87.90%

In this modality, InceptionV3 delivers the best performance, achieving the highest average accuracy and excelling in most individual traits. Extraversion and Openness continue to be the easiest traits to identify from first impressions, while

traits like conscientiousness and neuroticism, which depend on one's approach to various life domains, remain challenging to judge from initial encounters.

- **Comparative analysis of the results**

Inconsistent results are obtained when traits are analysed using multimodal decision models, as some of the models can analyse all the modalities together, while some of them fail to do so. This analysis (table 7.24) shows that InceptionV3+BERT Architecture is biased towards the video modality, hence showing results like those of Video models, while EfficientResNet + BERT architecture shows results closer to that of text models. The only model which could fit to all the modalities the best was XceptionResNet +BERT, which showed almost same accuracy on all labels except Neuroticism, thus making it the toughest trait to identify when we can fully assess the first impression from a video, undertaking all its modalities.

**Table 7.24.** Performance of multimodal decision models

Model	Avg. Accuracy	Accuracy on Various Traits				
		Extraversion	Openness	Agreeableness	Neuroticism	Conscientiousness
<b>InceptionV3 + BERT</b>	89.05%	89.94%	89.56%	88.46%	88.78%	88.49%
<b>EfficientResNet + BERT</b>	88.76%	89.56%	89.09%	88.65%	88.08%	88.39%
<b>XceptionResNet + BERT</b>	92.12%	92.35%	92.24%	92.16%	91.31%	92.51%

Our study provides valuable insights into the effectiveness of multimodal models for predicting personality traits. The choice of data modality significantly impacts model performance, with multimodal fusion models showing superior accuracy. This suggests that combining textual, visual, and auditory information offers a holistic understanding of an individual's personality. The ease of predicting certain traits, such as Extraversion, from textual and visual data implies that linguistic and non-verbal cues play a pivotal role in personality assessment. On the other hand, traits like Openness, Conscientiousness, and Neuroticism may require a more comprehensive assessment, considering multiple modalities. The multimodal fusion architecture, XceptionResNet + BERT, stands out as a robust choice for personality trait prediction, with consistently high accuracy across all traits. However, careful consideration of the modality and task at hand is essential when choosing the most appropriate model architecture.

Our system, leveraging the First Impression dataset, introduces several advancements in the field of personality-driven employment:

- **Enhanced Accuracy in Personality Assessment:** Utilizing multimodal data analysis, the system provides a more accurate and holistic view of an individual's personality traits, improving the precision of job-role matching.
- **Dynamic Adaptability:** By integrating real-time labor market analytics, the system can adapt to evolving job market trends, ensuring continuous alignment with current employment needs.
- **Scalable and Comprehensive Analysis:** The architecture is designed to handle large datasets, allowing for a comprehensive analysis of personality traits across diverse job sectors.

### 7.4.3. Performance of EmoMBTI-NET

Here we present the results of our comprehensive evaluation, illuminating the performance of diverse language models across both emoji mapping and MBTI personality classification tasks. Through meticulous analysis of Rouge scores and F1 metrics, clear patterns of strengths and weaknesses emerged, offering valuable insights into the nuanced interpretation of textual expressions and the efficacy of different models in capturing the intricacies of human communication.

The rouge scores, often used as a measure of similarity between the expected and the model-generated outcomes, were utilized here to gauge the efficacy of each model in emoji mapping (Table 7.25).

**Table 7.25.:** Emoji mapping Rouge scores using LLMs

Model	Final Training Loss (Rouge Score)
<b>Flan-T5</b>	<b>44.1280</b>
<b>Pegasus</b>	<b>11.5288</b>
<b>BART</b>	<b>0.0250</b>

The rouge scores reveal several key points about the performance of these models:

- **Higher Efficiency in Learning Representations:** BART's exceptionally low rouge score of 0.0250 suggests a superior ability to capture and reproduce the nuances required for emoji mapping. This implies that BART's training process effectively minimized the loss, leading to more accurate emoji predictions.
- **Comparison with Other Models:** Flan-T5, despite being a powerful and versatile model, showed a much higher final training loss in this task. This indicates that while Flan-T5 is generally effective across various tasks, it may not be as optimized as BART for the specific nuances of emoji mapping. Pegasus, known for its summarization capabilities, performed better than Flan-T5 but still fell short compared to BART. Its middle-range score suggests it has moderate capabilities in emoji mapping, potentially due to its underlying architecture which may prioritize certain aspects of language understanding differently.

- **Choice for MBTI Classification:** Given BART's outstanding performance in minimizing training loss, it was the logical choice for implementing the final emoji mappings for MBTI personality classification. The low rouge score indicates a high level of precision in BART's emoji mappings, which is crucial for accurately capturing personality traits through emojis.

In the task of classifying MBTI personalities using emojis, the accuracy of the models varied across different personality traits. Overall, BART consistently outperformed NLI-Roberta and NLI-DeBERTa, achieving the highest overall accuracy of 86.160% as shown in table 7.26.

**Table 7.26.:** Accuracy of different models

<i>Accuracy</i>	<b>NLI-Roberta</b>	<b>NLI-DeBERTa</b>	<b>BART</b>
<b>Overall</b>	0.78571	0.80654	0.8616
<b>Introversion</b>	0.84824	0.85504	0.88497
<b>Extroversion</b>	0.63736	0.63636	0.76331
<b>Intuition</b>	0.92096	0.92359	0.95000
<b>Sensing</b>	0.14035	0.18181	0.19318
<b>Thinking</b>	0.73437	0.76948	0.83742
<b>Feeling</b>	0.84308	0.83631	0.87172
<b>Judging</b>	0.73818	0.76351	0.81111
<b>Perceiving</b>	0.82897	0.83870	0.88356

Notably, BART displayed notable strengths in discerning traits related to introversion, intuition, and perceiving. However, the models demonstrated varying degrees of success in understanding traits such as extroversion and sensing.

The F1 scores presented in Table 7.27 for the models NLI-Roberta, NLI-DeBERTa, and BART provide an in-depth look at their performance across various MBTI personality traits. These scores are crucial for understanding how well each model balances precision and recall in their classifications, which is particularly important in personality analysis where both false positives and false negatives can lead to significant misinterpretations.

**Table 7.27.:** F1- scores of different models

<i>F-1</i>	<b>NLI-Roberta</b>	<b>NLI-DeBERTa</b>	<b>BART</b>
<b>Overall</b>	0.786	0.807	0.862
<b>Introversion</b>	0.848	0.855	0.885
<b>Extroversion</b>	0.637	0.636	0.763
<b>Intuition</b>	0.921	0.924	0.950
<b>Sensing</b>	0.140	0.182	0.193
<b>Thinking</b>	0.734	0.769	0.837
<b>Feeling</b>	0.843	0.836	0.872
<b>Judging</b>	0.738	0.764	0.811
<b>Perceiving</b>	0.829	0.839	0.884

Here's a detailed discussion of the F1 scores:

### Overall Performance

- **BART** shows the highest overall F1 score at **0.862**, indicating its superior general performance in classifying MBTI personalities.
- **NLI-DeBERTa** follows with **0.807**, and **NLI-Roberta** has an F1 of **0.786**. This suggests that while all models are relatively effective, BART excels in balancing precision and recall most effectively.

### Trait-Specific Performance

- **Introversion and Extroversion:**
  - BART outperforms in both traits, with particularly strong gains in Extroversion (**0.763**), which is notably higher than the scores for NLI-Roberta (**0.637**) and NLI-DeBERTa (**0.636**). This indicates BART's better capability in detecting features related to extroverted behavior through text.
- **Intuition:**
  - All models perform well with Intuition, with BART leading at **0.950**. This suggests that the models are effective at picking up on the abstract and theoretical content often associated with intuitive personalities.
- **Sensing:**
  - All models struggle with Sensing, with F1 scores remaining below **0.200**. This is likely due to the difficulty of capturing concrete and detail-oriented expressions through text, which are less prevalent or explicitly marked compared to intuitive expressions.
- **Thinking and Feeling:**
  - BART again leads in these traits, with Thinking at **0.837** and Feeling at **0.872**. These traits involve logical versus empathetic content, indicating BART's stronger ability to distinguish between structured argumentation and emotional expressions.
- **Judging and Perceiving:**
  - BART performs best in both Judging (**0.811**) and Perceiving (**0.884**), showing its effectiveness in identifying structured versus spontaneous expressions in text.

The results indicate BART's robustness in nuanced text interpretation, making it particularly suited for tasks requiring deep semantic understanding, such as personality assessment through text analysis. The superior F1 scores for BART

across most traits suggest that its training and underlying model architecture may be more attuned to the subtleties of personality-expressive language.

The lower scores in Sensing across all models highlight a common challenge in natural language processing capturing concrete and sensory-related information, which tends to be less explicitly expressed. This suggests a potential area for further model refinement and training, possibly by incorporating more sensory-specific training data or employing techniques that better capture these aspects. In conclusion, the F1 scores reveal the strengths and weaknesses of each model in handling different aspects of personality classification. BARTs leading performance across multiple traits recommends it as the preferable model for applications in personality analysis, offering more reliable and nuanced insights into individual personality types based on their textual expressions.

The choice of pre-trained language model and hyperparameters significantly influenced the performance of both emoji mapping and MBTI classification tasks. For emoji mapping, the hyperparameters such as learning rate and epochs were tuned for each LLM, optimizing their performance. Similarly, in MBTI classification, different learning rates and epochs were applied to NLI-Roberta, NLI-DeBERTa, and BART as shown in table 7.28.

**Table 7.28.:** Hyperparameters used

Model	Learning Rate	Epochs
<b>Emoji Mapping</b>		
Flan-T5	5e-5	3
Pegasus	3e-4	3
BART	1e-5	3
<b>MBTI Classification</b>		
NLI-Roberta	3e-5	10
NLI-DeBERTa	3e-6	10
BART	1e-5	10

Overall, the study highlights the effectiveness of large language models in leveraging emojis for personality detection tasks, with BART demonstrating superior performance in both emoji mapping and MBTI classification. The findings underscore the importance of model selection and hyperparameter tuning in optimizing the performance of such tasks and provide valuable insights into the nuances of interpreting personality traits conveyed through emojis.

### 7.5. Limitations of the study

Our exploration into the realm of automatic personality detection through deep learning models in user-generated content reveals several limitations that must be acknowledged to contextualize the findings and understand the potential constraints and future directions of this research.

### **7.5.1. Data Bias and Generalizability**

One of the primary limitations of our study is the inherent data bias present in the datasets utilized. A significant portion of these datasets are predominantly in English or are heavily skewed towards specific demographics, which introduces a notable bias. This linguistic and demographic skew raises substantial concerns regarding the generalizability of our findings to more diverse, global populations. For instance, cultural nuances and variations in personality expression across different languages and regions are not fully captured by the current datasets. Consequently, the models trained and tested on these datasets may not perform as accurately or effectively when applied to non-English speaking populations or culturally distinct groups. The over-representation of certain demographics can lead to models that are less sensitive to the unique personality expression patterns found in underrepresented groups, thereby limiting the applicability and fairness of our findings.

### **7.5.2. Interpretability Challenges**

Another significant limitation pertains to the interpretability of deep learning models. While these models exhibit powerful predictive capabilities, they often function as 'black boxes,' meaning that the internal decision-making processes are not transparent or easily understood. This lack of transparency is particularly problematic in the context of psychological applications, where understanding the reasoning behind model predictions is crucial. For instance, psychologists and other practitioners need to comprehend how and why a model arrived at a specific personality classification to ensure that the insights are valid and can be trusted. The opaque nature of deep learning models complicates this understanding, potentially undermining the trust and reliability of the models in real-world applications. Moreover, without clear interpretability, it is challenging to diagnose and correct potential biases or errors within the models, thereby limiting their overall efficacy and ethical application.

### **7.5.3. Data Quality and Annotation Consistency**

The quality and consistency of data annotation also present significant challenges. In personality detection tasks, accurate and consistent annotations are critical for training reliable models. However, the subjective nature of personality assessments can lead to variability in annotations, even among expert annotators. This variability can introduce noise and reduce the overall reliability of the training data, impacting the performance and robustness of the resulting models. Ensuring high-quality, consistent annotations across diverse datasets remains a complex and resource-intensive task that is crucial for the advancement of this field.

### **7.5.4. Scalability and Computational Resources**

Deep learning models, particularly those used for personality detection, require substantial computational resources for training and inference. This requirement can limit the scalability of the models and restrict their accessibility to researchers and practitioners with limited computational capabilities. The high computational



demands also pose challenges for deploying these models in real-time applications, where quick and efficient processing is essential. Developing methods to reduce the computational load while maintaining high performance is an ongoing challenge that needs to be addressed to enhance the practical applicability of these models.

#### **7.5.5. Ethical and Privacy Concerns**

Finally, ethical and privacy concerns are paramount in research involving personality detection from user-generated content. The use of personal data, especially sensitive information related to personality traits, necessitates stringent ethical standards and privacy protections. Ensuring that data collection, processing, and storage practices comply with ethical guidelines and legal regulations is essential to safeguard individuals' privacy and build trust in the technology. Furthermore, addressing potential biases in model predictions to prevent discriminatory outcomes is crucial for the ethical application of these models in various domains.

While our study has made significant strides in the field of automatic personality detection, these limitations highlight the need for ongoing efforts to address data bias, enhance model interpretability, ensure high-quality annotations, manage computational demands, and uphold ethical standards. Future research should focus on developing more inclusive and diverse datasets, improving the transparency of deep learning models, and implementing robust ethical frameworks to advance the responsible and effective use of personality detection technologies.

#### **7.6. Chapter Summary**

This chapter presented the results of evaluating various machine learning models for personality detection, such as PersonalityBERT and HindiPersonalityNet. These models were assessed using different linguistic datasets to determine their effectiveness in accurately identifying personality traits. These models were assessed using different linguistic datasets to determine their effectiveness in accurately identifying personality traits. Importantly, all models were thoroughly trained and validated ensuring the reliability of the performance metrics reported. The chapter also discussed the benefits of incorporating deep psychological theories into the models to enhance their predictive capabilities. It also addressed challenges encountered during the research, including data bias and computational demands, and summarized the outcomes as a foundation for future efforts to refine and expand the applicability of these models.

## CHAPTER 8

### CONCLUSION AND FUTURE WORK

Traversing the intricate landscape of personality research through computational analysis, this research has made several pivotal contributions. We developed and curated comprehensive datasets that capture personality traits across different languages and mediums, reflecting the multifaceted nature of personality expression. Our work has successfully introduced cross-linguistic models, which have demonstrated effective performance across various languages, showcasing their ability to predict and analyze personality traits in a culturally inclusive manner. Additionally, the research has illustrated practical applications, particularly in enhancing recruitment processes and organizational dynamics, where aligning personality traits with job requirements has led to actionable insights for improved hiring decisions and workplace efficiency.

Furthermore, the methodological enhancements achieved by integrating psychological insights with advanced computational techniques have resulted in more sophisticated and interpretable models, which significantly improve the precision and clarity of personality analyses. This interdisciplinary approach not only bridges the gap between psychology and technology but also paves the way for future innovations. Collectively, these contributions have expanded the knowledge base in both psychology and AI, offering new methodologies and frameworks that are both scientifically grounded and practically applicable. As we look to the future, these foundations will support continued exploration into the ethical and culturally sensitive expansion of personality detection technologies, ensuring they are robust, transparent, and accessible across diverse global contexts.

#### 8.1. Summary of Key Contributions

Throughout this research, several significant contributions have been made that advance the field of automatic personality detection. These contributions span from the creation of innovative datasets to the development of culturally adaptive models and practical applications that demonstrate the real-world impact of our work. Below are the primary contributions:

- **Creation of Comprehensive Datasets:** We developed and curated a diverse range of datasets capturing personality traits across different languages and mediums. These datasets include video clips, social media interactions, textual analysis, and emoji usage, reflecting the multifaceted nature of personality expression.
- **Development of Cross-Linguistic Models:** Our models have demonstrated effective performance across various languages, showcasing the ability to

predict and analyze personality traits in a culturally inclusive manner. These models have been validated for their accuracy and reliability, reinforcing their applicability in diverse contexts.

- **Practical Applications:** We have successfully illustrated real-world applications of our research, particularly in enhancing recruitment processes and organizational dynamics. By aligning personality traits with job requirements and team compositions, our models provide actionable insights that improve hiring decisions and workplace efficiency.
- **Methodological Enhancements:** Our interdisciplinary approach has integrated psychological insights with advanced computational techniques, resulting in more sophisticated and interpretable models. These improvements have enhanced the precision and clarity of our personality analyses.
- **Scientific Contributions:** Our work has made significant contributions to the fields of psychology and technology, offering new methodologies and frameworks that bridge these disciplines. These contributions pave the way for future research and technological innovations that incorporate psychological principles.

These contributions collectively advance the field of personality detection, offering robust, culturally adaptive models with significant real-world applications.

#### 8.1.1. Results of Research Objective 1

The first research objective focused on developing and validating models capable of accurately detecting personality traits across different languages and datasets. Below are the detailed outcomes for each model evaluated under this objective, highlighting their performance across various datasets.

- **Model: PersonalityBERT**
  - i. **Dataset:** English Kaggle\_MBTI
  - ii. **Performance:** PersonalityBERT achieves the highest accuracy among compared models, with a score of 69.45%. This outperforms both the LSTM model (38%) and the RNN model (67.77%) cited by Hernandez et al. The model's success underscores the efficacy of BERT-based approaches in processing nuanced language data for complex personality trait classification.
- **Model: ByaktitbaNet**
  - i. **Dataset:** Bangla dataset derived from "Baahubali 2"

- ii. **Performance:** ByaktitbaNet, combining BERT embeddings with LSTM, achieved an accuracy of 0.6849. This hybrid approach effectively classified dialogues into introvert, extrovert, and ambivert categories, demonstrating significant advancements in personality detection for low-resource languages.
- **Model: HindiPersonalityNet**
  - i. **Dataset:** Hindi Shaksiyat
  - ii. **Performance:** Utilizing BioWordVec embeddings with GRU, HindiPersonalityNet achieved an accuracy of 0.739. Compared to KBSVE-P's accuracy of 0.668 on a different Hindi MBTI dataset, HindiPersonalityNet shows superior performance, with higher F1-scores indicating balanced precision and recall.
- **Transformer-based Models for Personality\_Quotes**
  - i. **Dataset:** English Personality\_Quotes
  - ii. **Performance:** The highest accuracy of 0.791 was achieved by combining Common Crawl embeddings with the GRU model. Among transformer-based models, ELECTRA, DeBERTa, and BERT exhibited the highest performance with accuracy and F1-scores of approximately 0.8106. This suggests the robustness of these models in textual personality analysis.

### 8.1.2. Results of Research Objective 2

The second research objective aimed to enhance the performance of personality detection models by integrating advanced machine learning techniques and psychological theories. The following are the results of various models developed and tested, demonstrating their effectiveness in classifying personality traits across different datasets.

- **Model: KBSVE-P**
  - i. **Dataset:** Hindi Vishesh\_Charitr
  - ii. **Performance:** The KBSVE-P model, using various SVM kernels and voting techniques, demonstrated high accuracy in classifying MBTI personality types. The RBF kernel showed the best performance for N/S and T/F dimensions, while Linear and Polynomial kernels excelled in the I/E dimension. Soft Voting consistently yielded higher accuracy across most traits, highlighting its effectiveness.
- **Model: Decision Making MBTI**
  - i. **Dataset:** DM\_MBTI
  - ii. **Performance:** DeBERTa's superior performance in terms of accuracy (0.7395), precision (0.7426), recall (0.7395), F1-Score (0.74), and MCC score (0.6085), coupled with its balanced trade-off between

precision and recall, makes it the best choice for Decision Maker MBTI recognition.

- **ParvarishNET**
  - i. **Dataset:** Parvarish
  - ii. **Performance:** Tree-based machine learning algorithm, XGBoost demonstrates the highest performance than the Deep learning models, achieving an accuracy of 0.79 indicating the strength of the model in capturing complex patterns and interactions in data, leading to superior predictive performance.

### 8.1.3. Results of Research Objective 3

The third research objective focused on profiling emotional dispositions and integrating multimodal data to enhance personality detection. The results below highlight the performance of various models developed to analyze emotional attitudes and combine different data modalities, demonstrating their effectiveness in practical applications.

- **Model: Transformer-based for Manobhav Dataset**
  - i. **Dataset:** Hindi Manobhav Dataset
  - ii. **Performance:** Stacked mDeBERTa model achieved the highest F1 score of approximately 0.7745, indicating a strong balance between precision and recall. This model's success in profiling emotional attitudes like optimism and pessimism in Hindi underscores the benefits of stacking models and using ensemble techniques.
- **Model: Multimodal Fusion**
  - i. **Dataset:** Chalearn
  - ii. **Performance:** The fusion of XceptionResNet with BERT demonstrates exceptional performance, achieving a remarkable accuracy of 92.12% and other robust metrics. These empirical findings underscore the efficacy of the XceptionResNet + BERT approach in personality trait mapping, presenting an innovative and efficient method for job matching in urban environments.
- **EmoMBTI-NET**
  - i. **Dataset:** EmoMBTI
  - ii. **Performance:** BART's robustness in nuanced text interpretation, making it particularly suited for tasks requiring deep semantic understanding, such as personality assessment through text analysis. The superior F1 scores for BART across most traits suggest that its training and underlying model architecture may be more attuned to the subtleties of personality-expressive language.

Our exploration into automatic personality detection through deep learning models in user-generated content reveals several limitations, including data bias towards English and specific demographics, raising concerns about global generalizability, and interpretability challenges due to the 'black box' nature of deep learning models, which complicates understanding the decision-making process in psychological applications.

## **8.2. Social Impact**

Our research in personality detection has far-reaching social impacts across various domains. One of the most significant applications is in enhanced recruitment and team building. Businesses can leverage our models to identify personality traits that align with specific job roles, leading to better hiring decisions and more cohesive teams. For example, using personality detection tools, a company can match candidates who exhibit high conscientiousness and agreeableness for customer service roles, thereby enhancing workplace productivity and employee satisfaction. This targeted approach helps create harmonious work environments and reduces turnover rates.

In the realm of personalized marketing, our models enable marketing professionals to tailor their strategies to individual preferences, resulting in more effective and engaging campaigns. By understanding a consumer's personality, marketers can customize advertisements to resonate more deeply with the target audience. For instance, an extroverted individual might respond better to dynamic and social-oriented advertising, while an introverted person may prefer more informative and subtle marketing approaches. This personalization improves customer engagement and loyalty, ultimately driving better business outcomes.

Our research also has profound implications for mental health and well-being. The development of tools that accurately assess personality traits can aid mental health professionals in crafting personalized therapeutic approaches. For example, by identifying traits such as high neuroticism or low extraversion, therapists can tailor interventions to address specific emotional and behavioural patterns. Early identification of such traits can also facilitate the early detection of psychological issues, allowing for timely and effective interventions that enhance patient outcomes.

In educational settings, personality detection models can provide valuable insights into student behaviours and learning preferences. Educators can use these insights to personalize teaching methods, thereby improving educational outcomes. For example, understanding that a student with high openness to experience might thrive in a creative and exploratory learning environment allows teachers to tailor their instructional strategies accordingly. This personalized approach can foster a more engaging and effective educational experience for students.

Moreover, our work promotes cultural and linguistic inclusivity by developing models and datasets that cater to multiple languages and cultural contexts. This inclusivity ensures that the benefits of personality detection technology are accessible to diverse populations, addressing the need for culturally sensitive and linguistically appropriate tools. For instance, creating personality models for Hindi

and Bangla speakers helps extend psychological assessment and interventions to populations that have traditionally been underserved by predominantly English-based research. By embracing this diversity, our research contributes to a more equitable and comprehensive understanding of personality across the globe.

### **8.3. Recommendations for Future Research**

As we look to the future, several key areas of focus will guide our continued research and development:

- **Enhancing Model Interpretability**  
Enhancing the transparency and interpretability of our models is a primary goal. We aim to develop tools and techniques that allow users to understand the decision-making processes of our models, increasing trust and usability.
- **Reducing Technological Barriers**  
To democratize access to our tools, we plan to optimize our models to function effectively on less powerful hardware. This will make our research accessible to a wider audience, including those with limited computational resources.
- **Focus on Ethics and Privacy**  
Ensuring the ethical handling of personal information is paramount. We are committed to implementing robust privacy measures and ethical guidelines to protect the data we use and the individuals it represents, complying with data protection regulations and fostering transparency.
- **Expand Cross-Cultural Research**  
We intend to extend our research to include more languages and cultural contexts, further enhancing the global applicability of our models. This will involve curating new datasets from underrepresented regions and ensuring that our models are culturally sensitive and adaptable.
- **Integrate Multimodal Data**  
Future work will explore the integration of multimodal data sources, combining textual analysis with visual and auditory cues from video clips. This approach will provide a richer and more comprehensive understanding of personality traits, capturing the full spectrum of human expression.

In conclusion, our research has made significant strides in the field of personality analysis, developing new datasets, creating cross-linguistic models, and demonstrating practical applications. While recognizing the limitations of our current methodologies, we remain committed to advancing this field through continuous innovation and responsible research practices. Our future efforts will focus on enhancing model interpretability, reducing technological barriers, expanding cross-cultural research, and upholding the highest ethical standards. Through these



initiatives, we strive to develop models that are as interpretable as they are predictive, ultimately contributing to a deeper understanding of personality and its myriad expressions across different cultures and contexts.

#### **8.4. Closing Thoughts**

In conclusion, our journey through the intricate landscape of personality detection has yielded substantial advancements and insights, bridging the gap between psychology and technology. We have developed robust models and datasets, explored the depths of linguistic diversity, and demonstrated real-world applications that can significantly impact various sectors.

Our contributions underscore the potential of integrating psychological principles with advanced computational methods to create powerful tools for personality analysis. These tools not only enhance our scientific understanding of personality traits but also offer practical solutions for improving recruitment, marketing, mental health, education, and more.

Looking ahead, we remain committed to addressing the limitations identified in our study, such as data bias and interpretability challenges. Our future work will focus on making our models more transparent, reducing their computational requirements, and ensuring the ethical and responsible use of personal data. By continuing to innovate and refine our approaches, we aim to further advance the field of personality detection and contribute to the broader goal of enhancing human well-being through technology.

The intersection of personality research and computational analysis holds immense promise, and we are excited to be at the forefront of this transformative journey. Through continued collaboration, interdisciplinary research, and a steadfast commitment to ethical practices, we envision a future where personality detection technologies are not only highly accurate and reliable but also accessible and beneficial to all.



## References

1. Eysenck, H. J. (2012). *A Model for Personality*. New York: Springer Science & Business Media.
2. Wang, X., Zhao, Y., & Pourpanah, F. (2020). Recent advances in deep learning. *International Journal of Machine Learning and Cybernetics*, 11, 747-750.
3. Abu-Jbara, A., Ezra, J., & Radev, D. (2013, June). Purpose and polarity of citation: Towards nlp-based bibliometrics. In *Proceedings of the 2013 conference of the North American chapter of the association for computational linguistics: Human language technologies* (pp. 596-606).
4. Garousi, V., Bauer, S., & Felderer, M. (2020). NLP-assisted software testing: A systematic mapping of the literature. *Information and Software Technology*, 126, 106321.
5. Acheampong, F. A., Nunoo-Mensah, H., & Chen, W. (2021). Transformer models for text-based emotion detection: a review of BERT-based approaches. *Artificial Intelligence Review*, 54(8), 5789-5829.
6. Liu, X., Zheng, Y., Du, Z., Ding, M., Qian, Y., Yang, Z., & Tang, J. (2023). GPT understands, too. *AI Open*.
7. Shah, M. S., Bhat, M. A., Singh, M. S., Chavan, M. A., & Singh, M. A. (2010). Sentiment analysis.
8. Tyagi, P., & Tripathi, R. C. (2019, February). A review towards the sentiment analysis techniques for the analysis of twitter data. In *Proceedings of 2nd international conference on advanced computing and software engineering (ICACSE)*.
9. Anees, A. F., Shaikh, A., Shaikh, A., & Shaikh, S. (2020). Survey paper on sentiment analysis: Techniques and challenges. *EasyChair* 2516-2314.
10. Asghar, M. Z., Khan, A., Ahmad, S., & Kundi, F. M. (2014). A review of feature extraction in sentiment analysis. *Journal of Basic and Applied Scientific Research*, 4(3), 181-186.
11. Kumar, A., & Garg, G. (2020). Systematic literature review on context-based sentiment analysis in social multimedia. *Multimedia tools and Applications*, 79(21), 15349-15380.
12. Acharya, A., Singh, B., & Onoe, N. (2023, September). Llm based generation of item-description for recommendation system. In *Proceedings of the 17th ACM Conference on Recommender Systems* (pp. 1204-1207).
13. Sarker, I. H. (2024). LLM potentiality and awareness: a position paper from the perspective of trustworthy and responsible AI modeling. *Discover Artificial Intelligence*, 4(1), 40.
14. Keh, S.S., & Cheng, I. (2019). Myers-Briggs Personality Classification and Personality-Specific Language Generation Using Pre-trained Language Models. *ArXiv*, abs/1907.06333.
15. Ravichandiran, S. (2021). *Getting Started with Google BERT: Build and train state-of-the-art natural language processing models using BERT*. Packt Publishing Ltd.
16. Choong EJ, Varathan KD. **2021**. Predicting judging-perceiving of Myers-Briggs

Type Indicator (MBTI) in online social forum. *PeerJ* 9:e11382  
<https://doi.org/10.7717/peerj.11382>

17. Sonmezoz, K., Ugur, O., & Diri, B. (2020, October). MBTI personality prediction with machine learning. In 2020 28th Signal Processing and Communications Applications Conference (SIU) (pp. 1- 4). IEEE.
18. Amirhosseini, M. H., & Kazemian, H. (2020). Machine learning approach to personality type prediction based on the myers briggs type indicator®. *Multimodal Technologies and Interaction*, 4(1), 9.
19. Ismail, Sarerusaenye & Bashari Rad, Babak & Ismail, Shahrinaz. (2017). Significant of MBTI personality model on decision making in university program selection. 10.1109/ICITISEE.2017.8285560.
20. Cerkez, N., Vrdoljak, B., & Skansi, S. (2021). A method for MBTI classification based on impact of class components. *IEEE access*, 9, 146550-146567.
21. Vásquez, R. L., & Ochoa-Luna, J. (2021, October). Transformer-based approaches for personality detection using the MBTI model. In *2021 XLVII Latin American computing conference (CLEI)* (pp. 1-7). IEEE.
22. Behaz, A., & Djoudi, M. (2012). Adaptation of learning resources based on the MBTI theory of psychological types. *International Journal Of Computer Science Issues (IJCSI)*, 9(1), 135.
23. Boyle, G. J. (1995). Myers-Briggs type indicator (MBTI): some psychometric limitations. *Australian Psychologist*, 30(1), 71-74.
24. Jirásek, I., Janošíková, T., Sochor, F., & Česka, D. (2021). Some specifics of Czech recreation and leisure studies' students: Personality types based on MBTI. *Journal of Hospitality, Leisure, Sport & Tourism Education*, 29, 100315.
25. Cohen, Y., Ornoy, H., & Keren, B. (2013). MBTI personality types of project managers and their success: A field survey. *Project Management Journal*, 44(3), 78-87.
26. Gupta, N., Madhavan, A., Duvvuri, D., & Angeline, R. (2019). MBTI based personality prediction of a user based on their writing on social media. *Int J Eng Adv Technol (IJEAT)*.
27. Justindhas, Y., Mohanraj, S. M., & Shivani, R. (2022, August). A synoptic survey on personality prediction system using mbti. In *2022 3rd International Conference on Electronics and Sustainable Communication Systems (ICESC)* (pp. 950-955). IEEE.
28. Mushtaq, Z., Ashraf, S., & Sabahat, N. (2020, November). Predicting MBTI personality type with K-means clustering and gradient boosting. In *2020 IEEE 23rd International Multitopic Conference (INMIC)* (pp. 1-5). IEEE.
29. Gavrilescu, M., & Vizireanu, N. (2018). Predicting the Big Five personality traits from handwriting. *EURASIP Journal on Image and Video Processing*, 2018, 1-17.
30. Rammstedt, B., Roemer, L., Mutschler, J., & Lechner, C. (2023). The Big Five personality dimensions in large-scale surveys: An overview of 25 German data sets for personality research. *Personality Science*, 4(1), e10769.
31. Ning, H., Dhelim, S., & Aung, N. (2019). PersoNet: Friend recommendation system based on big-five personality traits and hybrid filtering. *IEEE Transactions on Computational Social Systems*, 6(3), 394-402.

32. Azucar, D., Marengo, D., & Settanni, M. (2018). Predicting the Big 5 personality traits from digital footprints on social media: A meta-analysis. *Personality and individual differences*, 124, 150-159.
33. Arya, R., Kumar, A., Bhushan, M., & Samant, P. (2022). Big five personality traits prediction using brain signals. *International Journal of Fuzzy System Applications (IJFSA)*, 11(2), 1-10.
34. Leutner, F., Ahmetoglu, G., Akhtar, R., & Chamorro-Premuzic, T. (2014). The relationship between the entrepreneurial personality and the Big Five personality traits. *Personality and individual differences*, 63, 58-63.
35. Gürpınar, F., Kaya, H., & Salah, A. A. (2016). Combining deep facial and ambient features for MBTIimpression estimation. In *Computer Vision–ECCV 2016 Workshops: Amsterdam, The Netherlands, October 8-10 and 15-16, 2016, Proceedings, Part III* 14 (pp. 372-385). Springer International Publishing.
36. Gürpınar, F., Kaya, H., & Salah, A. A. (2016, December). Multimodal fusion of audio, scene, and face features for first impression estimation. In *2016 23rd International conference on pattern recognition (ICPR)* (pp. 43-48). IEEE.
37. Gucluturk, Y., Guclu, U., Perez, M., Jair Escalante, H., Baro, X., Guyon, I., ... & Van Lier, R. (2017). Visualizing apparent personality analysis with deep residual networks. In *Proceedings of the IEEE International Conference on Computer Vision Workshops* (pp. 3101-3109).
38. Kaya, H., Gulpınar, F., & Ali Salah, A. (2017). Multi-modal score fusion and decision trees for explainable automatic job candidate screening from video cvs. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops* (pp. 1-9).
39. Junior, J. C. J., Güçlütürk, Y., Pérez, M., Güçlü, U., Andujar, C., Baró, X., ... & Escalera, S. (2019). First impressions: A survey on vision-based apparent personality trait analysis. *IEEE Transactions on Affective Computing*, 13(1), 75-95.
40. Yagmur Gucluturk, Umut Guclu, Xavier Baro, Hugo Jair Escalante, Isabelle Guyon, Sergio Escalera, Marcel A. J. van Gerven, and Rob van Lier. **2018**. Multimodal First Impression Analysis with Deep Residual Networks. *IEEE Trans. Affect. Comput.* 9, 3 (July 2018), 316\_329. <https://doi.org/10.1109/TAFFC.2017.2751469>.
41. Hassan, H. A. M., Sansonetti, G., Gasparetti, F., Micarelli, A., & Beel, J. (2019, September). Bert, elmo, use and infersent sentence encoders: The panacea for research-paper recommendation?. In *RecSys (Late-Breaking Results)* (pp. 6-10).
42. Alberti, C., Lee, K., & Collins, M. (2019). A bert baseline for the natural questions. *arXiv preprint arXiv:1901.08634*.
43. Acheampong, F. A., Nunoo-Mensah, H., & Chen, W. (2021). Transformer models for text-based emotion detection: a review of BERT-based approaches. *Artificial Intelligence Review*, 54(8), 5789-5829.
44. Kumar, A., Mallik, A., & Kumar, S. (2023). HumourHindiNet: Humour detection in Hindi web series using word embedding and convolutional neural network. *ACM Transactions on Asian and Low-Resource Language Information Processing*.
45. Katar, O., Özkan, D., Yıldırım, Ö., & Acharya, U. R. (2023). Evaluation of GPT-3 AI language model in research paper writing. *Turkish Journal of Science and Technology*, 18(2), 311-318.
46. Rai, N. (2016). *Bi-modal regression for Apparent Personality trait Recognition*.

- 2016 23rd International Conference on Pattern Recognition (ICPR). doi:10.1109/icpr.2016.7899607
47. Zhang, C. L., Zhang, H., Wei, X. S., & Wu, J. (2016, October). Deep bimodal regression for apparent personality analysis. In *European conference on computer vision* (pp. 311-324). Cham: Springer International Publishing.
  48. Hernández Y., Peña C.A., Martínez A. (2018) Model for Personality Detection Based on Text Analysis. In: Batyrshin I., Martínez-Villaseñor M., Ponce Espinosa H. (eds) *Advances in Computational Intelligence. MICAI 2018. Lecture Notes in Computer Science*, vol 11289. Springer, Cham. [https://doi.org/10.1007/978-3-030-04497-8\\_17](https://doi.org/10.1007/978-3-030-04497-8_17)
  49. Majumder, N., Poria, S., Gelbukh, A., & Cambria, E. (2017). Deep learning-based document modeling for personality detection from text. *IEEE Intelligent Systems*, 32(2), 74-79.
  50. Juan Abdon Miranda-Correa, Mojtaba Khomami Abadi, Nicu Sebe, and Ioannis Patras. 2021. *AMIGOS: A Dataset for Affect, Personality and Mood Research on Individuals and Groups*. *IEEE Trans. Affect. Comput.* 12, (April-June 2021), 479-493. <https://doi.org/10.1109/TAFFC.2018.2884461>.
  51. Kampman, Onno & Jebalbarez, Elham & Bertero, Dario & Fung, Pascale. (2018). Investigating Audio, Video, and Text Fusion Methods for End-to-End Automatic Personality Prediction. 606-611. 10.18653/v1/P18-2096.
  52. Sun, X., Liu, B., Meng, Q. *et al.* Group-level personality detection based on text generated networks. *World Wide Web* **23**, 1887–1906 (2020). <https://doi.org/10.1007/s11280-019-00729-2>
  53. Lynn, V., Balasubramanian, N., & Schwartz, H. A. (2020, July). Hierarchical modeling for user personality prediction: The role of message-level attention. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics* (pp. 5306-5316).
  54. Ren, Z., Shen, Q., Diao, X., & Xu, H. (2021). *A sentiment-aware deep learning approach for personality detection from text*. *Information Processing & Management*, 58(3), 102532. doi:10.1016/j.ipm.2021.102532
  55. Kamal El-Demerdash, Reda A. El-Khoribi, Mahmoud A. Ismail Shoman, Sherif Abdou, Psychological Human Traits Detection based on Universal Language Modeling, *Egyptian Informatics Journal*, Volume 22, Issue 3, 2021, Pages 239-244, ISSN 1110-8665, <https://doi.org/10.1016/j.eij.2020.09.001>.
  56. G.D. Salsabila, E.B. Setiawan, Semantic approach for big five personality prediction on twitter. *J. RESTI (Rekayasa Sistem Dan Teknologi Informasi)* **5**(4), 680–687 (2021)
  57. Singh, S., Singh, W. AI-based personality prediction for human well-being from text data: a systematic review. *Multimed Tools Appl* (2023). <https://doi.org/10.1007/s11042-023-17282-w>
  58. Ibrahim, R. T., & Ramo, F. M. (2023). Hybrid intelligent technique with deep learning to classify personality traits. *International Journal of Computing and*

- Digital Systems, 13(1), 231-244.
59. Sirasapalli, J.J., Malla, R.M. A deep learning approach to text-based personality prediction using multiple data sources mapping. *Neural Comput & Applic* 35, 20619–20630 (2023). <https://doi.org/10.1007/s00521-023-08846-w>
  60. Johnson, S.J., Murty, M.R. An Aspect-Aware Enhanced Psycholinguistic Knowledge Graph-Based Personality Detection Using Deep Learning. *SN COMPUT. SCI.* 4, 293 (2023). <https://doi.org/10.1007/s42979-023-01670-y>
  61. Yang, K., Lau, R. Y., & Abbasi, A. (2023). Getting personal: A deep learning artifact for text-based measurement of personality. *Information Systems Research*, 34(1), 194-222.
  62. Guo, A., Hirai, R., Ohashi, A. et al. Personality prediction from task-oriented and open-domain human–machine dialogues. *Sci Rep* 14, 3868 (2024). <https://doi.org/10.1038/s41598-024-53989-y>
  63. Grunenberg, E., Peters, H., Francis, M. J., Back, M. D., & Matz, S. C. (2024). Machine learning in recruiting: predicting personality from CVs and short text responses. *Frontiers in Social Psychology*, 1, 1290295.
  64. Sze, W. Y. S., Herrero, M. P., & Garriga, R. (2024). Personality Trait Inference Via Mobile Phone Sensors: A Machine Learning Approach. arXiv preprint arXiv:2401.10305.
  65. Suhartono, D., Ciputri, M. M., & Susilo, S. (2024). Machine Learning for Predicting Personality using Facebook-Based Posts. *Engineering, MAThematics and Computer Science Journal (EMACS)*, 6(1), 1-6.
  66. Hu, L., He, H., Wang, D., Zhao, Z., Shao, Y., & Nie, L. (2024, March). LLM vs Small Model? Large Language Model Based Text Augmentation Enhanced Personality Detection Model. In *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 38, No. 16, pp. 18234-18242).
  67. Serrano-Guerrero, J., Alshouha, B., Bani-Doumi, M., Chiclana, F., Romero, F. P., & Olivas, J. A. (2024). Combining machine learning algorithms for personality trait prediction. *Egyptian Informatics Journal*, 25, 100439.
  68. Saeidi, S. (2024). Identifying personality traits of WhatsApp users based on frequently used emojis using deep learning. *Multimedia Tools and Applications*, 83(5), 13873-13886.
  69. R. Liao, S. Song and H. Gunes, "An Open-source Benchmark of Deep Learning Models for Audio-visual Apparent and Self-reported Personality Recognition," in *IEEE Transactions on Affective Computing*, doi: 10.1109/TAFFC.2024.3363710.
  70. Alshouha, B., Serrano-Guerrero, J., Chiclana, F., Romero, F. P., & Olivas, J. A. (2024). Personality trait detection via transfer learning. *Comput Mater Continua*.
  71. Singh JK, Misra G, De Raad B. Personality structure in the trait lexicon of Hindi, a major language spoken in India. *European Journal of Personality*. 2013 Nov;27(6):605-20.



72. Singh JK, De Raad B. The personality trait structure in Hindi replicated. *International Journal of Personality Psychology*. 2017 Jun 29;3:26-35.
73. Khan SN, Leekha M, Shukla J, Shah RR. Vyaktiv: A multimodal peer-to-peer hindi conversations-based dataset for personality assessment. In 2020 IEEE Sixth International Conference on Multimedia Big Data (BigMM) 2020 Sep 24 (pp. 103-111). IEEE.
74. Khan, I. (2013). Big Five Personality Measurement Instrument-An Urdu Translation. Big Five Personality Measurement Instrument-An Urdu Translation.
75. Maghsoudi, M. (2013). Investigating the Effect of Big Five Personality Traits in Iranian EFL Bilingual Learners. *International Journal of Language and Linguistics*, 1(1), 26. <https://doi.org/10.11648/J.IJLL.S.20130101.15>
76. Adi GY, Tandio MH, Ong V, Suhartono D. Optimization for automatic personality recognition on Twitter in Bahasa Indonesia. *Procedia Computer Science*. 2018 Jan 1;135:473-80.
77. Yılmaz, T., Ergil, A., & İlgen, B. (2020). Deep learning-based document modeling for personality detection from Turkish texts. In *Proceedings of the Future Technologies Conference (FTC) 2019: Volume 1* (pp. 729-736). Springer International Publishing.
78. Rudra U, Chy AN, Seddiqui MH. Personality traits detection in bangla: A benchmark dataset with comparative performance analysis of state-of-the-art methods. In 2020 23rd International Conference on Computer and Information Technology (ICCIT) 2020 Dec 19 (pp. 1-6). IEEE.
79. Anari MS, Rezaee K, Ahmadi A. TraitLWNet: a novel predictor of personality trait by analyzing Persian handwriting based on lightweight deep convolutional neural network. *Multimedia Tools and Applications*. 2022 Mar;81(8):10673-93.
80. García-Peñalvo, F., Cruz-Benito, J., Martín-González, M., Vázquez-Ingelmo, A., Sánchez-Prieto, J. C., & Therón, R. (2018). Proposing a machine learning approach to analyze and predict employment and its factors.
81. Kumar, D., Verma, C., Singh, P. K., Raboaca, M. S., Felseghi, R. A., & Ghafoor, K. Z. (2021). Computational statistics and machine learning techniques for effective decision making on student's employment for real-time. *Mathematics*, 9(11), 1166.
82. Garg, S., Sinha, S., Kar, A. K., & Mani, M. (2022). A review of machine learning applications in human resource management. *International Journal of Productivity and Performance Management*, 71(5), 1590-1610.
83. Tiwari, R. (2023). The impact of AI and machine learning on job displacement and employment opportunities. *Interantional Journal of Scientific Research in Engineering and Management*, 7(01).
84. Wei, Y., Rao, X., Fu, Y., Song, L., Chen, H., & Li, J. (2023). Machine learning prediction model based on enhanced bat algorithm and support vector machine for slow employment prediction. *Plos one*, 18(11), e0294114.
85. Stimpson, A. J., & Cummings, M. L. (2014). Assessing intervention timing in computer-based education using machine learning algorithms. *IEEE Access*, 2, 78-87.

86. Halde, R. R. (2016, September). Application of Machine Learning algorithms for betterment in education system. In *2016 International Conference on Automatic Control and Dynamic Optimization Techniques (ICACDOT)* (pp. 1110-1114). IEEE.
87. Kučak, D., Juričić, V., & Đambić, G. (2018). MACHINE LEARNING IN EDUCATION-A SURVEY OF CURRENT RESEARCH TRENDS. *Annals of DAAAM & Proceedings*, 29.
88. Alenezi, H. S., & Faisal, M. H. (2020). Utilizing crowdsourcing and machine learning in education: Literature review. *Education and Information Technologies*, 25(4), 2971-2986.
89. Luan, H., & Tsai, C. C. (2021). A review of using machine learning approaches for precision education. *Educational Technology & Society*, 24(1), 250-266.
90. Munir, H., Vogel, B., & Jacobsson, A. (2022). Artificial intelligence and machine learning approaches in digital education: A systematic revision. *Information*, 13(4), 203.
91. Cho, G., Yim, J., Choi, Y., Ko, J., & Lee, S. H. (2019). Review of machine learning algorithms for diagnosing mental illness. *Psychiatry investigation*, 16(4), 262.
92. Abd Rahman, R., Omar, K., Noah, S. A. M., Danuri, M. S. N. M., & Al-Garadi, M. A. (2020). Application of machine learning methods in mental health detection: a systematic review. *Ieee Access*, 8, 183952-183964.
93. Tate, A. E., McCabe, R. C., Larsson, H., Lundström, S., Lichtenstein, P., & Kuja-Halkola, R. (2020). Predicting mental health problems in adolescence using machine learning techniques. *PloS one*, 15(4), e0230389.
94. Kim, J., Lee, D., & Park, E. (2021). Machine learning for mental health in social media: bibliometric study. *Journal of Medical Internet Research*, 23(3), e24870.
95. Chung, J., & Teo, J. (2022). Mental health prediction using machine learning: taxonomy, applications, and challenges. *Applied Computational Intelligence and Soft Computing*, 2022(1), 9970363.
96. Bai, Q., Dan, Q., Mu, Z., & Yang, M. (2019). A systematic review of emoji: Current research and future perspectives. *Frontiers in psychology*, 10, 476737.
97. Manganari, E. E. (2021). Emoji use in computer-mediated communication. *The International Technology Management Review*, 10(1), 1-11.
98. Kone, V. S., Anagal, A. M., Anegundi, S., Jadekar, P., & Patil, P. (2023). Emoji Prediction Using Bi-Directional LSTM. In *ITM Web of Conferences* (Vol. 53, p. 02004). EDP Sciences.
99. Rathod, J., Neha, K., Purohit, H. S., Verma, J., & Hiremath, S. (2023, July). Emoji Recommendation System Using Deep Learning Algorithms. In *2023 IEEE International Conference on Electronics, Computing and Communication Technologies (CONECCT)* (pp. 1-5). IEEE.
100. Philips, N., Sioen, I., Michels, N., Sleddens, E., & De Henauw, S. (2014). The influence of parenting style on health related behavior of children: findings from the ChiBS study. *International Journal of Behavioral Nutrition and Physical Activity*, 11, 1-14.
101. Sarwar, S. (2016). Influence of parenting style on children's behaviour. *Journal of Education and Educational Development*, 3(2).

102. Febiyanti, A., & Rachmawati, Y. (2021, March). Is authoritative parenting the best parenting style?. In *5th International Conference on Early Childhood Education (ICECE 2020)* (pp. 94-99). Atlantis Press.
103. El-Sawy, A., El-Bakry, H., & Loey, M. (2017). CNN for handwritten arabic digits recognition based on LeNet-5. In *Proceedings of the International Conference on Advanced Intelligent Systems and Informatics 2016 2* (pp. 566-575). Springer International Publishing.
104. Hossain, A., Konok, U. H., Islam, R., Ruhani, R. M. K., Musfikin, R., Uddin, M. M., ... & Tuhin, R. A. (2023, June). Utilizing GloVe Embeddings for Deep Learning-Based Analysis of Research Paper Abstracts. In *2023 5th International Congress on Human-Computer Interaction, Optimization and Robotic Applications (HORA)* (pp. 1-6). IEEE.
105. Dey, R., & Salem, F. M. (2017, August). Gate-variants of gated recurrent unit (GRU) neural networks. In *2017 IEEE 60th international midwest symposium on circuits and systems (MWSCAS)* (pp. 1597-1600). IEEE.
106. Tan, M., Santos, C. D., Xiang, B., & Zhou, B. (2015). Lstm-based deep learning models for non-factoid answer selection. *arXiv preprint arXiv:1511.04108*.
107. Siami-Namini, S., Tavakoli, N., & Namin, A. S. (2019, December). The performance of LSTM and BiLSTM in forecasting time series. In *2019 IEEE International conference on big data (Big Data)* (pp. 3285-3292). IEEE.
108. Bolukbasi, T., Chang, K. W., Zou, J., Saligrama, V., & Kalai, A. (2016). Quantifying and reducing stereotypes in word embeddings. *arXiv preprint arXiv:1606.06121*.
109. Basaldella, M., & Collier, N. (2019, November). Bioreddit: Word embeddings for user-generated biomedical NLP. In *Proceedings of the Tenth International Workshop on Health Text Mining and Information Analysis (LOUHI 2019)* (pp. 34-38).
110. Panchenko, A., Ruppert, E., Faralli, S., Ponzetto, S. P., & Biemann, C. (2017). Building a web-scale dependency-parsed corpus from CommonCrawl. *arXiv preprint arXiv:1710.01779*.
111. Yamada, I., Asai, A., Sakuma, J., Shindo, H., Takeda, H., Takefuji, Y., & Matsumoto, Y. (2018). Wikipedia2Vec: An efficient toolkit for learning and visualizing the embeddings of words and entities from Wikipedia. *arXiv preprint arXiv:1812.06280*.
112. Stevens, J., Chen, D., Zimmer, J., Punturo, B., & Kim, M. (2019). Representing document-level semantics of biomedical literature using pre-trained embedding models. In *Proceedings of the 15th Conference on Natural Language Processing, KONVENS*.
113. Goeckenjan, G., Sitter, H., Thomas, M., Branscheid, D., Flentje, M., Griesinger, F., ... & Deppermann, K. (2011). PubMed results. *Pneumologie*, 65(8), e51-e75.
114. Shahbaz, M., Suresh, L., Rexford, J., Feamster, N., Rottenstreich, O., & Hira, M. (2019). Elmo: Source routed multicast for public clouds. In *Proceedings of the ACM Special Interest Group on Data Communication* (pp. 458-471).



115. Yao, T., Zhai, Z., & Gao, B. (2020, March). Text classification model based on fasttext. In *2020 IEEE International conference on artificial intelligence and information systems (ICAIS)* (pp. 154-157). IEEE.
116. Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R. R., & Le, Q. V. (2019). Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems*, 32.
117. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... & Stoyanov, V. (2019). Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
118. Clark, K., Luong, M. T., Le, Q. V., & Manning, C. D. (2020). Electra: Pre-training text encoders as discriminators rather than generators. *arXiv preprint arXiv:2003.10555*.
119. He, P., Liu, X., Gao, J., & Chen, W. (2020). Deberta: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*.
120. Matthews, G., & Gilliland, K. (1999). The personality theories of HJ Eysenck and JA Gray: A comparative review. *Personality and Individual differences*, 26(4), 583-626.
121. Buhmann, M. D. (2000). Radial basis functions. *Acta numerica*, 9, 1-38.
122. Aizawa, A. (2003). An information-theoretic perspective of tf-idf measures. *Information Processing & Management*, 39(1), 45-65.
123. Zhang, Y., Zhang, H., Cai, J., & Yang, B. (2014). A weighted voting classifier based on differential evolution. In *Abstract and applied analysis* (Vol. 2014, No. 1, p. 376950). Hindawi Publishing Corporation.
124. Pires, T., Schlinger, E., & Garrette, D. (2019). How multilingual is multilingual BERT?. *arXiv preprint arXiv:1906.01502*.
125. Ta, H. T., Rahman, A. B. S., Najjar, L., & Gelbukh, A. F. (2022). Transfer Learning from Multilingual DeBERTa for Sexism Identification. In *IberLEF@SEPLN*.
126. Mohsin, M. A., & Beltiukov, A. (2019, May). Summarizing emotions from text using Plutchik's wheel of emotions. In *7th Scientific Conference on Information Technologies for Intelligent Decision Making Support (ITIDS 2019)* (pp. 291-294). Atlantis Press.
127. Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P., & Soricut, R. (2019). Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*.
128. Farag, H. H., Said, L. A., Rizk, M. R., & Ahmed, M. A. E. (2021). Hyperparameters optimization for ResNet and Xception in the purpose of diagnosing COVID-19. *Journal of Intelligent & Fuzzy Systems*, 41(2), 3555-3571.
129. Xia, X., Xu, C., & Nan, B. (2017, June). Inception-v3 for flower classification. In *2017 2nd international conference on image, vision and computing (ICIVC)* (pp. 783-787). IEEE.
130. Lu, S., Hong, Q., Wang, B., & Wang, H. (2020). Efficient resnet model to predict protein-protein interactions with gpu computing. *IEEE Access*, 8, 127834-127844.

131. Sanh, V., Debut, L., Chaumond, J., & Wolf, T. (2019). DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.
132. Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., ... & Zettlemoyer, L. (2019). Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*.

## APPENDIX A: LIST OF PUBLICATIONS WITH PROOF

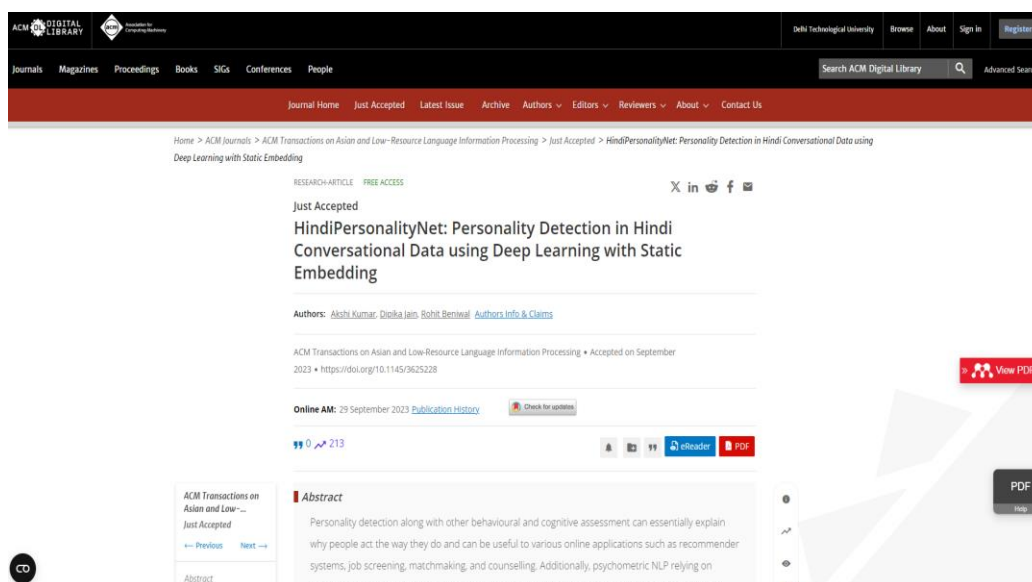
### SCIE JOURNALS

[1] Kumar, A., Beniwal, R, Jain, D. (2023). Personality Detection using Kernel-based Ensemble Model for leveraging Social Psychology in Online Networks *ACM Transactions on Asian and Low-Resource Language Information Processing (ACM TALLIP)*- <https://doi.org/10.1145/3571584> [SCIE-Impact Factor: 1.472]

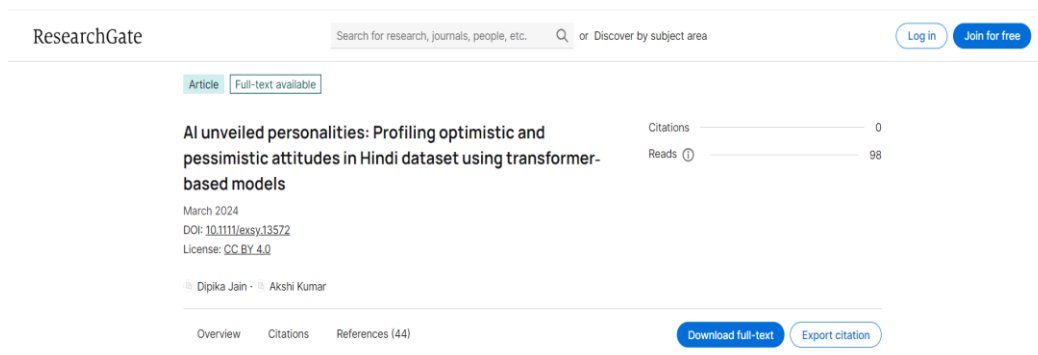
PUBLISHED

The screenshot shows the ACM Digital Library interface. At the top, there's a navigation bar with links like 'Journals', 'Magazines', 'Proceedings', 'Books', 'SIGs', 'Conferences', and 'People'. Below this is a search bar and a 'Search ACM Digital Library' button. The main content area displays the title 'Personality Detection using Kernel-based Ensemble Model for Leveraging Social Psychology in Online Networks' by authors Akshi Kumar, Rohit Beniwal, and Diptika Jain. It includes the journal information 'ACM Transactions on Asian and Low-Resource Language Information Processing, Volume 22, Issue 5' and the article number 'Article No.: 151, pp 1-20'. The publication date is '09 May 2023'. There are links for 'Check for updates', 'eReader', and 'PDF'. The abstract is visible, starting with 'The Asian social networking market dominates the world landscape with the highest consumer penetration rate...'. On the right side, there's a 'View PDF' button and a 'PDF' button. The bottom left corner shows a 'CD' logo.

[2] Kumar, A., Jain, D., & Beniwal, R. (2023). HindiPersonalityNet: Personality Detection in Hindi Conversational Data using Deep Learning with Static Embedding. *ACM Transactions on Asian and Low-Resource Language Information Processing*. <https://doi.org/10.1145/3625228> [SCIE-Impact Factor: 1.472] **PUBLISHED**



[3] Jain, D., & Kumar, A., (2024). AI unveiled personalities: Profiling optimistic and pessimistic attitudes in Hindi dataset using transformer-based models, *Expert Systems*, John Wiley & Sons <https://doi.org/10.1111/exsy.13572> [SCIE-Impact Factor: 2.812], **PUBLISHED**



## UGC-CARE JOURNAL

1. Jain, D., Beniwal, R. & Kumar, A. (2024). Advancements in Personality Detection: Unleashing the Power of Transformer-Based Models and Deep Learning with Static Embeddings on English Personality Quotes. *International Journal of All Research Education & Scientific Methods, (IJARESM)*, <https://doi.org/10.56025/IJARESM.2023.1201242235> [UGC-Care]

## PUBLISHED



**International Journal of All Research  
Education & Scientific Methods**  
An ISO Certified Peer-Reviewed Journal

ISSN: 2455-6211  
Convert Your Language ▾



HOME EDITORIAL BOARD PROCESSING CHARGES ONLINE SUBMISSION ISSUES INDEXING CONTACT US

Search here 

**IJARESM Menu**

- Publication Ethics
- Peer Review & Publication Policy
- Call For Papers
- Why IJARESM
- Topics Covered
- Special Issue

**Download**

- Author Guidelines
- Copyrights Form
- Paper Template

## Search Result

You Are Here :  > Search Result

Total Records : 18 Records ▾

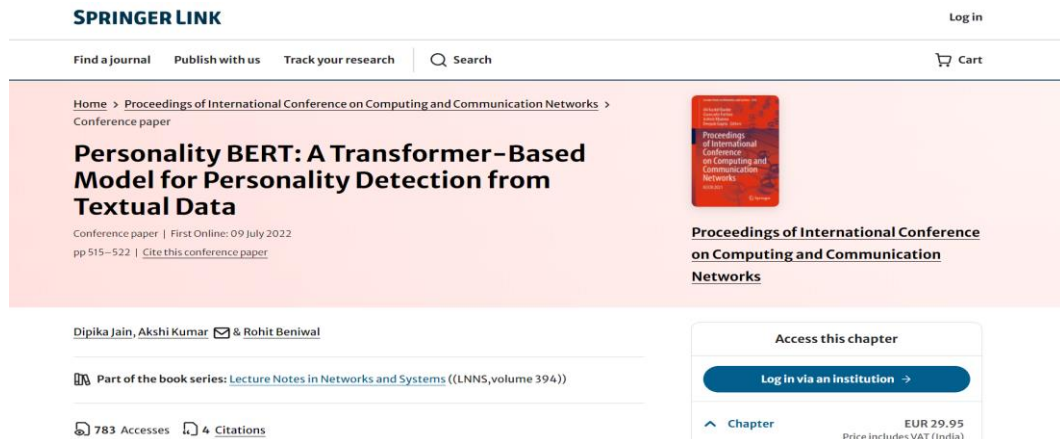
<b>Title</b>	Advancements in Personality Detection: Unleashing the Power of Transformer-Based Models and Deep Learning with Static Embeddings...
<b>Author</b>	Dipika Jain, Dr. Rohit Beniwal, Dr. Akshi Kumar
<b>Country</b>	India

[View](#) [DOWNLOAD](#)

Total Records : 18 Records ▾

## CONFERENCES

1. Jain, D., Kumar, A., Beniwal, R. (2021). Personality BERT: A Transformer-Based Model for Personality Detection from Textual Data. *In Proceedings of International Conference on Computing and Communication Networks. Lecture Notes in Networks and Systems*, vol 394. Springer, Singapore.  
[https://doi.org/10.1007/978-981-19-0604-6\\_48](https://doi.org/10.1007/978-981-19-0604-6_48)



2. Jain, D., Beniwal, R., Kumar, A. (2023). ByaktitbaNet: Deep Neural Network for Personality Detection in Bengali Conversational Data. *In Proceedings of Fourth Doctoral Symposium on Computational Intelligence. DoSCI 2023. Lecture Notes in Networks and Systems*, vol 726. Springer, Singapore.  
[https://doi.org/10.1007/978-981-99-3716-5\\_57](https://doi.org/10.1007/978-981-99-3716-5_57)



## APPENDIX B: PLAGIARISM REPORT

### Similarity Report

PAPER NAME

**final chapters thesis.docx**

WORD COUNT

**35642 Words**

CHARACTER COUNT

**220390 Characters**

PAGE COUNT

**104 Pages**

FILE SIZE

**4.8MB**

SUBMISSION DATE

**Aug 20, 2024 10:01 AM GMT+5:30**

REPORT DATE

**Aug 20, 2024 10:02 AM GMT+5:30**

#### ● 9% Overall Similarity

The combined total of all matches, including overlapping sources, for each database.

- 5% Internet database
- 4% Publications database
- Crossref database
- Crossref Posted Content database
- 6% Submitted Works database

#### ● Excluded from Similarity Report

- Bibliographic material
- Quoted material
- Cited material
- Small Matches (Less than 8 words)
- Manually excluded text blocks