# Who Are You Behind the Screen?
# Implicit MBTI and Gender Detection Using Artificial Intelligence

Kourosh Shahnazari[1*†] and Seyed Moein Ayyoubzadeh[1†]

[1*]Computer Engineering Department, Sharif University of Technology.

*Corresponding author(s). E-mail(s): kourosh@null.net;
Contributing authors: smoein.ayyoubzadeh16@sharif.edu;
[†]These authors contributed equally to this work.

## Abstract

In personalized technology and psychological research, precisely detecting demographic features and personality traits from digital interactions becomes ever more important. This work investigates implicit categorization, inferring personality and gender variables directly from linguistic patterns in Telegram conversation data, while conventional personality prediction techniques mostly depend on explicitly self-reported labels. We refine a Transformer-based language model (RoBERTa) to capture complex linguistic cues indicative of personality traits and gender differences using a dataset comprising 138,866 messages from 1,602 users annotated with MBTI types and 195,016 messages from 2,598 users annotated with gender. Confidence levels help to greatly raise model accuracy to 86.16%, hence proving RoBERTa's capacity to consistently identify implicit personality types from conversational text data. Our results highlight the usefulness of Transformer topologies for implicit personality and gender classification, hence stressing their efficiency and stressing important trade-offs between accuracy and coverage in realistic conversational environments. With regard to gender classification, the model obtained an accuracy of 74.4%, therefore capturing gender-specific language patterns. Personality dimension analysis showed that people with introverted and intuitive preferences are especially more active in text-based interactions. This study emphasizes practical issues in balancing accuracy and data coverage as Transformer-based models show their efficiency in implicit personality and gender prediction tasks from conversational texts.

1

# 1 Introduction

## 1.1 The Myers-Briggs Type Indicator (MBTI)

The Myers-Briggs Type Indicator (MBTI) is a widely utilized psychometric instrument in psychological research, therapy, corporate management, and educational settings for categorizing personality types based on individual preferences. The MBTI, further developed and systematized by Katharine Cook Briggs and her daughter Isabel Briggs Myers in the mid-20th century, stemmed from the psychological types proposed by Carl Gustav Jung in 1921. Analyzing natural tendencies across several psychological dimensions enables the characterization of human behavior, cognition, and social interactions based on foundational theoretical frameworks.[1]

The MBTI evaluates personal inclinations across four opposing dimensions:

- **Extraversion (E) - Introversion (I)** This dimension distinguishes an individual's focus on the exterior environment from their internal domain. Extraverted individuals typically succeed in social interactions, exhibiting active engagement with others and their environment. Introverted individuals gain energy and comfort from reflection, solitary activities, and reflective meditation, often requiring serene environments for productivity and personal growth.
- **Sensing (S) – Intuition (N)** This dimension contrasts an individual's preference for tangible, empirical, and observable information (Sensing) with abstract, conceptual, and theoretical insights (Intuition). Individuals who prefer sensation rely heavily on direct sensory experiences and practical knowledge. Conversely, persons who choose intuition typically engage with concepts, patterns, and potential situations, often prioritizing innovation and theoretical frameworks over current realities.
- **Thinking (T) – Feeling (F)**: This component reflects how individuals participate in decision-making processes. Individuals inclined towards Thinking prioritize logical, objective analysis and consistent principles, emphasizing rationality and impartiality. Individuals who select Feeling are guided by empathy, personal values, and the quest for emotional balance, emphasizing subjective experiences and social relationships.
- **Judging (J) - Perceiving (P)**: This dimension indicates individuals' perspectives on structure, organization, and flexibility in their daily lives. Individuals who prioritize Judging typically favor structured, regulated, and organized environments, often exhibiting decisiveness and a desire for closure. In contrast, individuals who prefer Perceiving demonstrate flexibility, spontaneity, and receptiveness, often thriving in uncertainty and improvisation.

The synthesis of these four opposing attributes yields sixteen distinct personality types, represented by combinations such as INTJ, ESFP, INFP, or ESTJ. Each kind represents distinct behavioral patterns, cognitive styles, interpersonal dynamics, and professional or academic inclinations [1]. Thus, accurately identifying and understanding an individual's MBTI type provides significant insights into their cognitive processes, preferences, and interpersonal dynamics, substantially improving psychological evaluation, career guidance, team leadership, and educational practices.

## 1.2 Transformer-Based Language Models

In the last ten years, substantial advancements have occurred in Natural Language Processing (NLP), primarily due to the emergence of Transformer-based neural network architectures. The Transformer, introduced by Vaswani et al. in 2017 [2], has revolutionized language modeling by replacing traditional recurrent neural networks (RNNs) and convolutional neural networks (CNNs) with a framework exclusively based on self-attention mechanisms. This self-attention allows the model to dynamically evaluate the importance of different words inside a sentence or context, enhancing the modeling of semantic relationships and context-dependent linguistic features.

Three notable Transformer-based models have set new standards for several NLP tasks, including text classification, sentiment analysis, and language generation:

- **BERT (Bidirectional Encoder Representations from Transformers)**: Introduced by Devlin et al. in 2018 [3], BERT profoundly transformed NLP approaches by capturing contextual information bidirectionally. Through the utilization of masked language modeling and next-sentence prediction tasks during pre-training on large datasets, BERT attains an enhanced comprehension of context, semantics, and syntactic subtleties, leading to its widespread adoption in many NLP tasks, including question answering and text summarization.
- **RoBERTa (Robustly Optimized BERT Approach)**: RoBERTa, developed by Liu et al. in 2019 [4], builds upon BERT by incorporating optimizations such as dynamic masking during training, an expanded training dataset, and the elimination of the next-sentence prediction target. These improvements enable RoBERTa to apprehend nuanced linguistic subtleties and contextual variability, consistently surpassing BERT and other baseline models across many NLP benchmarks.
- **GPT-2 (Generative Pre-trained Transformer 2)**: GPT-2, introduced by Radford et al. [5], is fundamentally distinct from BERT and RoBERTa due to its unidirectional architecture, which is primarily focused on generating coherent and contextually pertinent text. GPT-2 has unparalleled capabilities in language generation, encompassing text completion and creative writing endeavors. Notwithstanding its generative emphasis, GPT-2 may be adeptly fine-tuned for discriminative classification tasks, illustrating the versatility of Transformer-based models.

The efficacy of Transformer models in natural language processing has significant ramifications, especially in utilizing linguistic analysis for personality and demographic profiling. By precisely capturing intricate linguistic patterns, these models offer a potential toolkit for assessing textual data to deduce psychological traits.

## 1.3 Implicit Data Collection and Its Significance

The significant proliferation of digital communication platforms such as social media, forums, and messaging applications has generated vast repositories of textual content created either overtly or implicitly by users. Analyzing this passive data stream has profound implications for psychological research, as it enables the covert extraction of personal characteristics, such as personality types and gender, without explicit self-reporting.

This method of implicit data collection presents numerous advantages:

- **Enhanced User Profiling**: Linguistic analysis improves user profiles and provides personality insights that improve service personalizing including customized learning platforms and targeted marketing strategies [6].
- **Mental Health Monitoring**: Automated language analysis of user interactions helps to early detect and observe psychological disorders including depression or anxiety, so providing important opportunities for quick intervention and support [7].
- **Ethical and Privacy Considerations**: Although implicit data collecting has great possibilities, it also raises important ethical questions about user permission, privacy protection, and data governance. Thus, careful thought and open conversation on data use are absolutely crucial [8].[6].

## 1.4 Research Objectives

Situated at the junction of computational linguistics and psychological assessment, this study has as its main goals:

- Evaluating the effectiveness of Transformer-based language models (BERT, RoBERTa, and GPT-2) for accurately classifying MBTI personality types and gender using textual data.
- Investigating and identifying specific linguistic features and patterns that significantly contribute to accurate personality and gender classification.
- Discussing the broader societal, ethical, and practical implications of employing sophisticated NLP models for implicit psychological profiling and passive demographic identification.

Intending to greatly advance both NLP and psychological science and so foster an interdisciplinary dialogue vital for responsible innovation and practical application in psychological profiling and personalized technology, this research analyzes extensive datasets including 138,666 messages from 1,602 users for MBTI classification and 195,016 messages from 2,598 users for gender classification.

# 2 Related Work

## 2.1 Personality Prediction and MBTI Classification in Text

Automated personality classification is a crucial subject in computational linguistics and applied natural language processing (NLP) for the psychological and social sciences. The Myers-Briggs Type Indicator (MBTI) is a prevalent psychological framework that classifies individuals into 16 unique personality types, derived from the intersections of four dichotomies: Introversion (I) vs. Extraversion (E), Sensing (S) vs. Intuition (N), Thinking (T) vs. Feeling (F), and Judging (J) vs. Perceiving (P). This framework renders MBTI a commonly utilized instrument for linguistic and behavioral analysis [9, 10]. Recent improvements in transformer designs have prompted studies to focus on utilizing deep learning models like BERT, RoBERTa, ALBERT, and DistilBERT for personality detection in textual data.

## 2.2 Transformer-Based MBTI Classification

Utilizing Transformer-based models for MBTI categorization has demonstrated significant advancements compared to conventional machine learning and lexicon-based approaches [11]. Previous models predominantly utilized feature-based methodologies, including TF-IDF vectorization, LIWC dictionaries, and manually generated linguistic indicators. These approaches possessed intrinsic limitations, as they did not adequately capture contextual representations or semantic links within textual data.

Recent improvements indicate that pre-trained models, especially BERT-based architectures, exhibit superior performance in retrieving personality-related linguistic cues. Vásquez and Ochoa Luna [11] investigated Transformer-based MBTI categorization utilizing the Kaggle dataset, attaining exemplary results with 88.63% accuracy and 88.97% F1-score. Their research established that contextual embeddings surpass conventional word-vector models, emphasizing BERT's capacity to discern subtle personality-related expressions across MBTI categories.

Research by does Santos and Paraboni [12] shown that fine-tuning BERT for MBTI recognition on social media datasets markedly surpassed traditional classifiers, attaining strong generalization across several test settings. Kadambi [13] validated these findings in their research on Twitter users with self-identified MBTI types, demonstrating that user profiles, status updates, and liked tweets uniquely enhance personality categorization accuracy.

## 2.3 Limitations of Explicit MBTI Classification

Most MBTI classification research employ datasets in which individuals explicitly self-report their personality type, resulting in intrinsic biases and concerns over data validity. Tareaf [14] and Julianda and Maharani [15] examined personality classification models developed using Reddit data, noting considerable class imbalances in MBTI distributions—certain personality types (e.g., ISTJ, ENFP) were markedly more prevalent than others, resulting in overfitting on majority-class samples.

To resolve this issue, Li et al. [10] developed MBTIBench, a meticulously maintained dataset employing soft-labeling techniques to mitigate inaccuracies arising from self-reported misclassifications. Their research revealed that approximately 30% of data samples had erroneous user-assigned MBTI types, highlighting a significant issue with the dependability of self-reported personality datasets.

Notwithstanding enhancements from dataset curation methodologies like MBTIBench, this research continues to depend on explicit personality labels, which fail to translate effectively to implicit personality classification tasks in real-world scenarios, such as chat-based conversations.

## 2.4 Implicit Personality Prediction in Conversational Text

Unlike prior MBTI classification efforts, our work focuses on implicit personality prediction, where personality types are inferred without explicit self-reported labels. This presents multiple challenges:

- **Absence of Direct Labels**: In contrast to Reddit and Twitter datasets, which feature user-provided MBTI annotations, Telegram data does not contain clear MBTI type disclosures. Consequently, personality assessment must deduce characteristics solely from linguistic styles, conversational involvement, and behavioral tendencies.
- **Noise and Informal Linguistic Structures**: Transforming natural conversation into effective model input involves substantial preprocessing, as Telegram messages contain slang, abbreviations, emojis, and multilingual text [12, 16].
- **Overlapping Personality and Gender Markers**: Several studies highlight correlations between MBTI dimensions and gender-linked linguistic traits, requiring careful disentangling of these features to avoid overfitting [16].

## 2.5 Transformer Fine-Tuning for Conversational Data

Recent research highlights the need for robust fine-tuning procedures when applying Transformers to noisy data sources:

- Arya et al. [17] conducted comparative analyses on BERT, RoBERTa, DistilBERT, and ALBERT for personality classification, demonstrating that BERT-based models consistently produce the highest accuracy across social media text.
- Applications of hybrid filtering frameworks, including dynamic tokenization pipelines and structured pre-processing for informal digital conversations, are crucial for adapting models to real-world messaging platforms [15].
- Evaluation metrics must handle class imbalances dynamically. Prior studies apply macro-F1 and weighted precision-recall scores to counteract dataset bias. Methods such as oversampling, class reweighting, and contrastive learning have been explored as potential solutions [9, 10].

## The Need for Implicit Personality Modeling

The prevalence of explicit self-reported MBTI classifications in previous studies considerably restricts the application of these models in authentic conversational contexts. Our research addresses this issue by (1) eliminating self-labeled training data, (2) enhancing tokenization and preprocessing methods for chat-based corpora, and (3) implementing Transformer-based fine-tuning techniques for implicit personality detection in Telegram conversations. This research transcends prior MBTI classification frameworks, facilitating the advancement of more resilient, real-time personality inference in text-based digital communication.

# 3 Methodology

# Transformer Models: Mechanism and Evolution

The emergence of **Transformer** models transformed natural language processing (NLP) by implementing a very efficient architecture founded on the self-attention mechanism. In contrast to conventional *Recurrent Neural Networks (RNNs)* and *Long Short-Term Memory (LSTM) networks*, Transformers eradicate sequential constraints

by utilizing *parallelized training* throughout whole input sequences. The basic design proposed by Vaswani et al. [2] employs a multi-layered self-attention mechanism that adaptively modifies weight distributions to encapsulate contextual word dependencies.

## 3.1 Self-Attention and Multi-Head Attention

In Transformer models, each input token contributes to the overall representation based on its *attention scores* across all other tokens in a sequence. The attention score is computed using the following key operations:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \tag{1}$$

where $Q$ (Query), $K$ (Key), and $V$ (Value) matrices represent different transformations of input representations, and $d_k$ is the scaling factor. **Multi-head attention** extends this mechanism by maintaining multiple attention weight matrices, allowing the model to capture diverse linguistic relationships.

## 3.2 Pretraining Strategies: Masked Language Modeling and Next Sentence Prediction

Most modern Transformer models rely on *unsupervised pretraining techniques*, enabling them to learn extensive general-purpose linguistic features before fine-tuning on specific tasks. Notable pretrained architectures include:

- **BERT (Bidirectional Encoder Representations from Transformers)** [3]: Introduces **masked language modeling (MLM)**, where random tokens in input text are masked, forcing the model to predict masked words using surrounding context.
- **RoBERTa (Robustly Optimized BERT Approach)** [4]: Enhances BERT's training regime by removing next-sentence prediction (NSP) and training on significantly larger datasets.
- **GPT (Generative Pretrained Transformer)** [5]: Utilizes a causal autoregressive decoding mechanism, allowing real-time text generation. Unlike BERT, GPT processes input unidirectionally.

## 3.3 Expanding Transformer Applications in NLP

Transformers have significantly advanced several NLP applications, including:

- **Text Classification**: Models like BERT, RoBERTa, and DistilBERT dominate various classification tasks, from *sentiment analysis* to *document categorization* [4].
- **Machine Translation**: The Transformer model became the backbone of leading *neural machine translation (NMT) systems*, outperforming LSTMs in real-time sentence translation [2].
- **Conversational AI**: Large-scale Transformer models, such as *GPT-based architectures*, power chatbots, virtual assistants, and real-time dialogue agents [5].
- **Psychological and Social Insights**: Transformers contribute to personality detection, author profiling, and *mental health prediction via linguistic cues* [6, 7].

7

## 3.4 Why Transformers Are Ideal for Personality and Gender Classification

Several Transformer properties make them particularly well-suited for personality classification from text:

1. **Bidirectional Context Awareness**: Unlike RNNs and LSTMs, BERT-based architectures process entire text sequences bidirectionally, capturing *long-range dependencies* in user discourse.
2. **Self-Attention for Linguistic Variability**: Personality and gender classification tasks benefit from *self-attention mechanisms*, allowing models to highlight influential grammatical and stylistic cues.
3. **Scalability and Adaptability**: Transformers can efficiently fine-tune on personality-labeled datasets while adapting to domain-specific text, such as *Telegram chat-based data* [11, 15].

## 3.5 The Need for Domain-Specific Transformer Adaptation

While prior research has demonstrated the efficacy of Transformers in psychological linguistics, most studies use data from structured or explicitly labeled user-generated content (e.g., *Reddit or Twitter MBTI-labeled datasets* [12, 13]). However, **Telegram conversations introduce unique linguistic challenges**:

- **Multi-Turn Conversations and Threading** Unlike Reddit, where posts exist independently, Telegram chat contexts evolve dynamically. Standard Transformer *tokenization strategies struggle* to segment multi-turn dialogues effectively.
- **Noise and Non-Standard Grammar** Informal messages include abbreviated words, emojis, and conversation artifacts that require *more sophisticated text filtering techniques* than standard document classification tasks.
- **Ethical and Privacy Challenges** Personality classification from chats raises ethical concerns regarding user profiling and data privacy [8]. Responsible data anonymization, differential privacy techniques, and user consent protocols should inform future work.

## Integrating Transformer Advances for Implicit Personality Detection

The Transformer revolution has enabled substantial progress in NLP-based personality classification; nonetheless, current research predominantly centers on explicit personality disclosures, such as self-reported MBTI classifications. Our research seeks to enhance the discipline by:

- Utilizing the contextual representations of Transformers for chat-based implicit personality inference.
- Enhancing preprocessing/tokenization techniques to accommodate informal Telegram communication.
- Guaranteeing comprehensive model evaluation through multi-layered fine-tuning strategies.

The integration of self-attention architectures, unsupervised personality modeling, and informal conversational adaptations will facilitate the development of more realistic psychological AI applications, connecting theoretical NLP advancements with practical behavioral insights in real-world dialogues.

# 4 Overview of the Transformer Architecture

The Transformer model follows an **encoder-decoder** structure:

1. The **encoder** processes input sequences into contextual representations.
2. The **decoder** generates output sequences by attending to encoder outputs and previous generated tokens.

Both components consist of multiple stacked layers, incorporating:

- **Multi-head self-attention mechanisms**
- **Feedforward networks**
- **Layer normalization and residual connections**

# 5 Step-by-Step Breakdown of the Transformer Algorithm

## 5.1 Input Embedding and Positional Encoding

Since Transformers **do not have recurrence**, they require a method to retain positional order in sequences. The model uses:

1. Learned **word embeddings** to convert input tokens into dense vectors.
2. **Positional encoding** to introduce order information, computed as follows:

$$PE_{\text{pos},2i} = \sin\left(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}\right) \tag{2}$$

$$PE_{\text{pos},2i+1} = \cos\left(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}\right) \tag{3}$$

where pos represents the position index and $d_{\text{model}}$ is the dimensionality of embeddings.

## 5.2 Encoder Mechanism

Each encoder layer performs the following:

### 5.2.1 Multi-Head Self-Attention

Each token representation is transformed into three matrices:

$$Q = XW_Q, \quad K = XW_K, \quad V = XW_V \tag{4}$$

where:

9

- $Q$ (Query) determines what information to extract.
- $K$ (Key) contains relevance scores against other tokens.
- $V$ (Value) represents token information passed through attention weightings.

The **scaled dot-product attention** is computed as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \tag{5}$$

where the denominator $\sqrt{d_k}$ prevents instability in gradient updates.

### 5.2.2 Residual Connections and Layer Normalization

To stabilize training, the model adds a **residual connection**:

$$Z = \text{LayerNorm}(A + X) \tag{6}$$

Following this, a **feedforward network** enhances non-linear transformations:

$$O = \text{ReLU}(ZW_1 + b_1)W_2 + b_2 \tag{7}$$

The final encoder output is normalized:

$$E_X = \text{LayerNorm}(O + Z) \tag{8}$$

## 5.3 Decoder Mechanism

The decoder consists of similar layers to the encoder but adds:

1. **Masked self-attention** to prevent information leakage from future tokens.
2. **Cross-attention with encoder outputs** to integrate input sequence information.

### 5.3.1 Masked Attention

To ensure autoregressive sequence generation, future tokens are masked:

$$A = \text{MaskedMultiHeadAttention}(Q, K, V) \tag{9}$$

### 5.3.2 Encoder-Decoder Cross-Attention

Each decoder token attends to encoder outputs:

$$A' = \text{MultiHeadAttention}(Q', K', V') \tag{10}$$

where:

- $Q'$ comes from the decoder.
- $K'$ and $V'$ come from the encoder outputs.

10

## 5.4 Final Output Computation

The decoder output is passed through a Softmax layer to compute word probabilities:

$$P(y_t) = \text{Softmax}(E_Y W_o) \tag{11}$$

The predicted sequence is obtained via:

$$\hat{Y} = \arg\max P(y_t) \tag{12}$$

---

**Algorithm 1** Transformer: Encoder-Decoder Architecture

---

**Require:** Input sequence $X = (x_1, x_2, ..., x_n)$, Target sequence $Y = (y_1, y_2, ..., y_m)$
**Ensure:** Transformed output sequence $\hat{Y}$

1: **Initialization:** Load parameters $\theta$ for Encoder and Decoder networks
2: **Compute input embeddings:** $E_X \leftarrow \text{Embed}(X) + \text{PosEnc}(X)$
3: **Compute target embeddings:** $E_Y \leftarrow \text{Embed}(Y) + \text{PosEnc}(Y)$
4: **for each encoder layer** $l = 1$ to $L$ **do**
5:     $Q, K, V \leftarrow E_X W_Q^l, E_X W_K^l, E_X W_V^l$            ▷ Linear projections
6:     $A \leftarrow \text{MultiHeadAttention}(Q, K, V)$            ▷ Self-Attention
7:     $Z \leftarrow \text{LayerNorm}(A + E_X)$     ▷ Residual connection and Layer Normalization
8:     $O \leftarrow \text{FeedForward}(Z)$            ▷ Position-wise feedforward network
9:     $E_X \leftarrow \text{LayerNorm}(O + Z)$            ▷ Final residual connection
10: **end for**
11: Encoder output: $H \leftarrow E_X$
12: **for each decoder layer** $l = 1$ to $L$ **do**
13:     $Q, K, V \leftarrow E_Y W_Q^l, E_Y W_K^l, E_Y W_V^l$
14:     $A \leftarrow \text{MaskedMultiHeadAttention}(Q, K, V)$ ▷ Masked to prevent seeing future tokens
15:     $Z \leftarrow \text{LayerNorm}(A + E_Y)$
16:     $Q' \leftarrow Z W_{Q'}^l, K' \leftarrow H W_{K'}^l, V' \leftarrow H W_{V'}^l$
17:     $A' \leftarrow \text{MultiHeadAttention}(Q', K', V')$ ▷ Cross-Attention with encoder output
18:     $Z' \leftarrow \text{LayerNorm}(A' + Z)$
19:     $O \leftarrow \text{FeedForward}(Z')$
20:     $E_Y \leftarrow \text{LayerNorm}(O + Z')$
21: **end for**
22: Decoder output: $\hat{Y} \leftarrow \text{Softmax}(E_Y W_o)$            ▷ Final token probabilities
23: **return** $\hat{Y}$

---

## 5.5 Dataset and Preprocessing

Our study utilizes an anonymized dataset of messages extracted from a Telegram social channel. Users have self-reported their MBTI personality type and gender. Unlike structured datasets, this dataset consists of naturally occurring conversations, allowing us to explore implicit personality and gender classification through linguistic features.
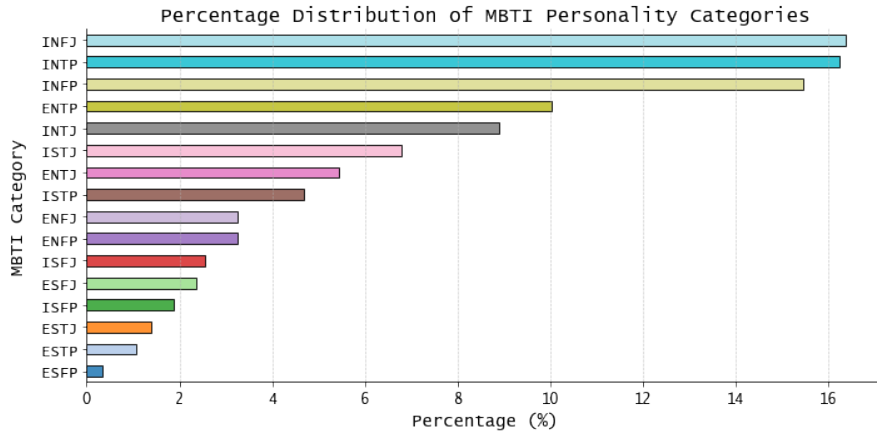
**Fig. 1** Distribution of messages across MBTI types and gender categories in the dataset.

## Dataset Composition

The dataset is divided into two primary tasks:

- **MBTI Classification**: 138,866 messages from 1,602 users with annotated MBTI labels.
- **Gender Classification**: 195,016 messages from 2,598 users with annotated gender labels.

Each message contains:

1. The raw message text.
2. An associated MBTI type (if available).
3. An associated gender label (if available).

Some messages include both MBTI and gender labels, while others contain only one of these attributes.

## MBTI Feature Distribution

We show a representation of the distribution of MBTI personality traits in our dataset, focusing on the four key dichotomies: Extroversion/Introversion, Sensing/Intuition, Thinking/Feeling, and Judging/Perceiving. Comprehending these distributions is essential.

The distribution of MBTI personality traits in our sample, illustrated in Figure 2, demonstrates significant variability across the four dichotomous dimensions.

These discrepancies may be shaped by the intrinsic preferences of persons possessing particular personality traits regarding computer-mediated communication (CMC), such as messaging platforms.
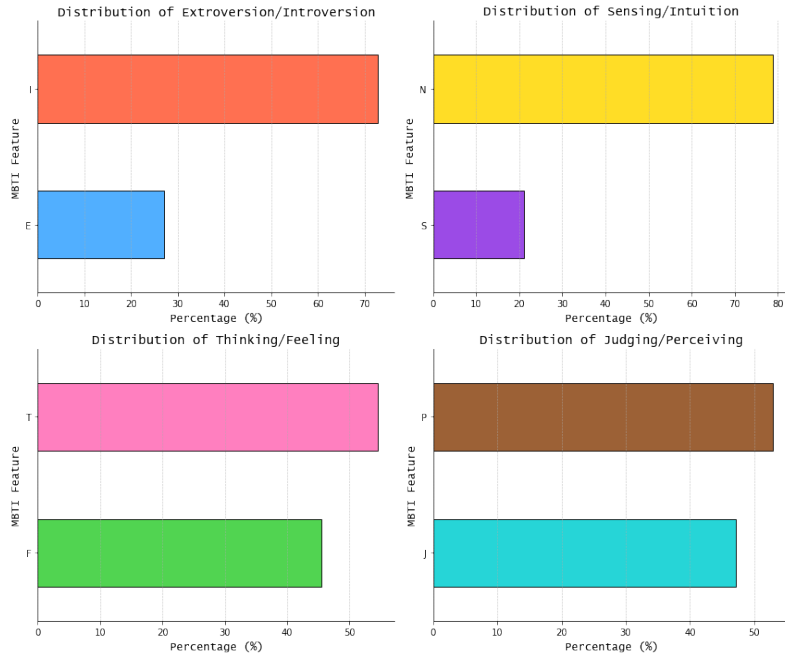
**Fig. 2** Distribution of MBTI personality dichotomies within the dataset. Each subplot represents the proportion of messages associated with a specific MBTI feature.

### 5.5.1 Introversion (I) Preference

Individuals exhibiting Introverted (I) inclinations frequently prefer textual communication to in-person contact. This desire enables individuals to regulate social interactions at their own tempo, along with their propensity for introspection and deliberate self-presentation. Studies suggest that introverts tend to favor online communication approaches, as these platforms allow for contemplative processing prior to responding [18]. As a result, introverts may exhibit more activity on messaging platforms, resulting in a greater representation of Introversion (I) within the dataset.

### 5.5.2 Intuition (N) Prevalence

Individuals with an intuitive preference are drawn to abstract concepts and future possibilities, often engaging in discussions that explore theoretical ideas and patterns. Digital communication platforms Online messaging applications are well-suited for investigations, as they enable extensive discussions at any time, allowing individuals to articulate their thoughts without immediate concerns regarding the practical implications of their statements. The discussion may result in a higher representation of Intuitive (N) users within our sample.

### 5.5.3 Thinking (T) Dominance

The higher prevalence of Thinking (T) relative to Feeling (F) in the sample can be attributed to the nature of online communication, which often emphasizes information sharing and discussion over emotional expression. Individuals who prioritize logic and objectivity may find messaging systems advantageous for participating in analytical and problem-solving discussions, thereby improving their engagement and representation.

### 5.5.4 Perceiving (P) Notability

Perceiving (P) individuals are characterized by adaptability and spontaneity, traits that align well with the dynamic and flexible nature of online messaging. The lack of rigid structure in CMC allows Perceivers (P) to navigate conversations fluidly, aligning with their preference for open-endedness and improvisation. This compatibility may contribute to the notable presence of Perceiving (P) users in the dataset.
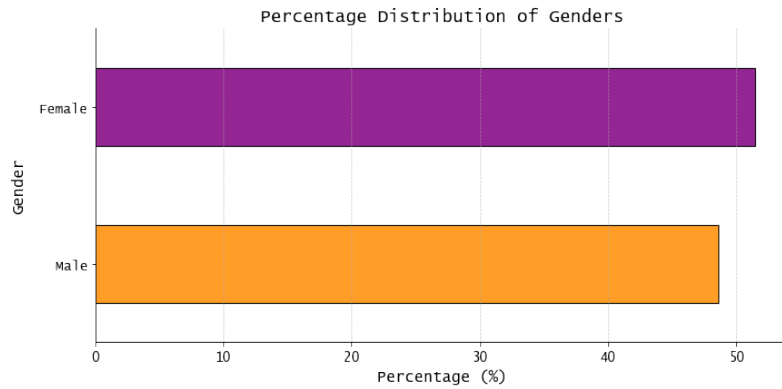


**Fig. 3** Percentage distribution of genders in the dataset. The chart represents the proportion of messages attributed to male and female users.

The dataset's distribution of MBTI personality types, illustrated in Figure 1, shows significant differences compared to general population distributions. INFJ is identified as the most common type, with INTP and INFP following closely, whereas ESFP exhibits the lowest representation. This pattern reveals distinct differences in communication styles driven by personality, self-selection biases, and interaction preferences in online messaging forums.

INFJs, representing roughly 16.37% of the dataset, exhibit a pronounced preference for text-based communication, attributed to their introspective characteristics and inclination towards profound, meaningful interactions. Their inclination towards structured and emotionally engaging discussions corresponds with the nature of online forums, enabling them to articulate their insights and values thoughtfully, free from immediate social pressures.

INTPs account for 16.22%, indicative of their analytical and intellectually curious characteristics. INTPs are drawn to platforms that provide intellectual stimulation,

excelling in settings that allow for the analysis of complex ideas and participation in theoretical discussions. However, their generally selective participation style may result in a reduced volume of interactions compared to INFJs, who tend to engage more frequently due to their relational approach.

INFPs comprise 15.44% of the dataset, reflecting a notable presence attributed to their inclination towards self-expression and the exploration of abstract concepts. Their introspective and emotionally expressive communication style effectively engages online forums through detailed and reflective posts. This corresponds with their established inclination to utilize written communication as a means of examining personal insights and emotions.

ENTP, at 10.02%, and INTJ, at 8.89%, exemplify personality types characterized by a focus on intellectual engagement and abstract reasoning. ENTPs exhibit a strong inclination towards debate and the exploration of ideas, often taking the initiative or engaging actively in dynamic discussions. INTJs favor conceptual and strategic discussions, utilizing forums as platforms for idea exchange, free from the distractions of in-person interactions.

ISTJ (6.79%) and ENTJ (5.43%) types exhibit moderate levels of engagement. The communication style of ISTJs is characterized by practicality and structure, which may restrict their engagement in abstract or emotionally charged discussions, as they tend to prioritize factual and clearly defined exchanges. ENTJs typically adopt a leadership-oriented approach that emphasizes strategic contributions over frequent participation, reflecting their goal-directed and succinct communication style.

ESFJ (2.37%) and ISFJ (2.55%) demonstrate lower representation, consistent with their inclination towards direct, socially structured interactions. Their focus on face-to-face and relational communication may hinder substantial involvement in online text-based discussions, which lack interpersonal cues.

The least represented types—ESTJ (1.40%), ESTP (1.08%), and ESFP (0.34%)—exhibit reduced activity levels, attributed to their pronounced inclination towards immediate, pragmatic, and socially interactive environments. These personality types emphasize direct interaction, active engagement, and sensory experiences, which diminishes their tendency to participate in written and reflective online discussions.

The MBTI type distribution in our dataset highlights the impact of personality traits on digital communication behaviors. Types defined by introspection, abstraction, and a preference for textual engagement (INFJ, INTP, INFP) are significantly overrepresented, whereas action-oriented, socially driven types (ESFP, ESTP, ESTJ) are notably underrepresented. The findings underscore significant personality-driven variations that affect online participation.

## Gender Distribution

The dataset's gender distribution, illustrated in Figure 3, presents the ratio of messages from male and female users. The dataset includes contributions from both genders, exhibiting no significant dominance of one category over the other. This balanced distribution allows the model to encounter a variety of linguistic patterns linked to different genders.

Linguistic research indicates that gender differences in language use can be observed through variations in sentence structure, word choice, and conversational style. These distinctions may function as valuable attributes for the gender classification task. The identification of gender-specific linguistic patterns within the dataset allows Transformer-based models to effectively discern features that differentiate male and female communication styles.

The dataset offers a comprehensive representation of both genders; however, it is crucial to recognize that language use is shaped by various factors beyond gender, such as personality, cultural background, and social context. Consequently, although gender-based linguistic tendencies may appear in model predictions, they should not be regarded as definitive indicators of gender identity. The classification model learns probabilistic relationships between textual features and the assigned gender labels.

The insights derived from this distribution enhance our understanding of the influence of gender representation in textual data on classification performance. Future research may examine biases in gender prediction models and assess the extent to which specific linguistic features disproportionately influence classification outcomes.

### 5.5.5 Preprocessing Pipeline

Given Telegram messages' informal and unstructured nature, substantial preprocessing was performed to standardize text and reduce noise. Our preprocessing pipeline includes:

1. **Removing Links**: All hyperlinks were removed to eliminate irrelevant content.
2. **Filtering Short Messages**: Messages with fewer than 25 tokens were discarded to ensure meaningful linguistic representation.
3. **Eliminating Explicit MBTI Mentions**: Any direct mention of MBTI types (e.g., "I am an ISTJ") was removed to prevent information leakage.
4. **Tokenization**: We applied model-specific tokenization methods, such as Word-Piece for BERT and Byte-Pair Encoding (BPE) for GPT-2.

These preprocessing steps ensure that only relevant linguistic features contribute to model training.

## 5.6 Model Selection and Feature Extraction

To classify MBTI types and gender, we employ Transformer-based models, extracting high-dimensional feature vectors from their embedding layers.

### 5.6.1 Transformer Models for Feature Extraction

We utilize the following pre-trained Transformer architectures:

- **BERT (Bidirectional Encoder Representations from Transformers)**
- **RoBERTa (Robustly Optimized BERT Pretraining Approach)**
- **GPT-2 (Generative Pretrained Transformer 2)**

For each model, we extract feature vectors from the final embedding layer, capturing contextual linguistic representations. The extracted feature vectors serve as input for classification.

### 5.6.2 Feature Vector Representation

The extracted embeddings are processed with the following configuration:

- **Hidden size**: 768
- **Number of hidden layers**: 12
- **Number of attention heads**: 12
- **Output vector dimensionality**: 16

These embeddings encapsulate complex linguistic structures and are used as input features for MBTI and gender classification.

## 5.7 Fine-Tuning Strategy

Each model undergoes fine-tuning using a classification head attached to the final embedding representation. The fine-tuning process follows these steps:

1. **Tokenization**: Messages are tokenized using the model's corresponding tokenizer.
2. **Feature Extraction**: Hidden-state embeddings from the final layer are used as classification inputs.
3. **Dimensionality Reduction**: The extracted 768-dimensional embeddings are reduced to a 16-dimensional representation.
4. **Classification Head**: A fully connected layer with a softmax activation function predicts the target label.

The loss function for fine-tuning is categorical cross-entropy, defined as:

$$L = -\sum_{i=1}^{N} y_i \log(\hat{y}_i) \tag{13}$$

where $y_i$ is the true label and $\hat{y}_i$ is the predicted probability.

### 5.7.1 Training Configuration

The training process is configured as follows:

- **Batch size**: 64
- **Learning rate**: $2.86e^{-5}$, optimized using AdamW
- **Number of epochs**: 16
- **Max length (Tokens)**: 128
- **Graphics Processing Unit (GPU)**: Nvidia Tesla T4

## 5.8 Evaluation Metrics

To assess classification performance, we utilize:

1. **Accuracy**: Measures the proportion of correctly classified instances.

2. **Precision, Recall, and F1-Score**: Evaluates class-specific performance, balancing false positives and false negatives.
3. **Macro-F1 Score**: Accounts for class imbalance by averaging F1-scores across all MBTI types.
4. **Confusion Matrix Analysis**: Identifies misclassification patterns and model biases.

# 6 Results

In this section, we present the performance evaluation of our models on MBTI personality classification and gender classification tasks. We report accuracy, precision, recall, and F1-score as key performance metrics. Table 1 and Table 4 summarize our findings.

## 6.1 MBTI Personality Classification

Table 1 presents the performance of BERT, GPT-2, and RoBERTa on the MBTI personality classification task. Among these models, RoBERTa demonstrates the highest classification performance, achieving an accuracy of **49%** and an F1-score of **50%**. This suggests that RoBERTa effectively captures the linguistic nuances associated with different MBTI personality types, outperforming both BERT and GPT-2.

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| BERT | 40% | 39% | 40% | 0.39 |
| GPT-2 | 45% | 45% | 45% | 0.44 |
| RoBERTa | **49%** | **51%** | **49%** | **0.50** |

**Table 1** Performance comparison for MBTI personality classification.

To further analyze the model's classification behavior, we visualize the confusion matrix for RoBERTa in Figure 4. The confusion matrix provides insight into misclassifications across MBTI types, revealing that certain personality categories, particularly those with overlapping linguistic patterns, are more prone to misclassification.

## 6.2 Effect of Confidence Thresholds on MBTI Classification

In standard classification scenarios, models assign the class with the highest probability to each instance, irrespective of the confidence level of the prediction. In practical applications, particularly those related to implicit personality classification, the reliability of predictions is essential. In response to this issue, we implemented a confidence threshold parameter, categorizing only instances with predicted confidence surpassing a defined threshold.

Modifying the confidence threshold serves as an effective method for markedly improving prediction accuracy. RoBERTa effectively captures robust linguistic indicators that are strongly associated with specific MBTI personality types by selectively

classifying instances with higher confidence scores. The results indicate that at a confidence threshold of 0.99, RoBERTa attained an accuracy of 86.16%. This significant enhancement highlights the model's efficacy in recognizing and utilizing linguistic patterns that are indicative of specific personality traits.

Table 2 presents the relationship between different confidence thresholds, classification accuracy, and the proportion of classified data. As shown, increasing the confidence threshold leads to a notable increase in classification accuracy. However, this comes at the cost of reduced data coverage. For example, at a threshold of 0.50, the model achieves 52.83% accuracy while covering 85% of the dataset, whereas at the highest threshold of 0.99, accuracy reaches 86.16% but covers only 26% of the dataset.

**Table 2** Effect of Confidence Thresholds on MBTI Classification

| Threshold | Accuracy (%) | Data Coverage (%) |
|---|---|---|
| 0.50 | 52.83 | 85 |
| 0.60 | 55.51 | 77 |
| 0.70 | 59.82 | 68 |
| 0.80 | 64.12 | 59 |
| 0.90 | 70.16 | 48 |
| 0.99 | 86.16 | 26 |

The high accuracy observed at elevated thresholds has significant implications, particularly in domains where precision is critical, such as personalized psychological counseling, targeted marketing, and educational settings. Employing elevated confidence thresholds improves reliability in automated personality predictions, guaranteeing that the outcomes or insights generated by the model are precise and applicable.

Adopting higher confidence thresholds leads to a reduction in the number of classified data points, highlighting a trade-off between accuracy and data coverage. As seen in Table 2, at the highest threshold of 0.99, approximately 26% of the total dataset remains classified. This suggests that although RoBERTa's predictions exhibit high accuracy at increased confidence levels, a significant amount of data remains unclassified.

Nevertheless, this limitation is offset by the immense value that such a highly accurate classification provides. The ability to predict an individual's MBTI personality type with 86.16% accuracy, solely based on their text messages, is a remarkable achievement in the field of computational psychology and AI-driven personality analysis. This level of precision suggests that even in real-world applications where only a fraction of messages can be confidently classified, the insights generated by the model remain highly reliable and valuable.

Such a high-confidence personality prediction system can be particularly useful in applications where accurate personality profiling is essential, such as in psychological assessments, AI-driven recruitment systems, and adaptive learning platforms that personalize educational content based on inferred personality traits. The ability to

infer personality traits implicitly, without requiring users to take explicit psychological tests, not only enhances user experience but also opens the door to more seamless and non-intrusive applications in personalized AI systems.

Furthermore, our results demonstrate the substantial benefits of confidence-threshold tuning in transformer-based implicit personality classification tasks, emphasizing the methodological rigor and practical relevance of our approach. By dynamically adjusting the confidence threshold, practitioners can optimize their models for either broad coverage or high precision, depending on their specific application requirements.

A significant avenue for future research involves utilizing multiple messages per user instead of focusing on single messages, which may enhance both prediction accuracy and coverage through aggregation. Considering multiple conversational samples per individual allows models to address the coverage limitation seen at high confidence thresholds, thereby improving the practical applicability and robustness of implicit personality detection from conversational data. By leveraging larger conversational contexts, it is likely that we can push classification accuracy even further while maintaining broader coverage.
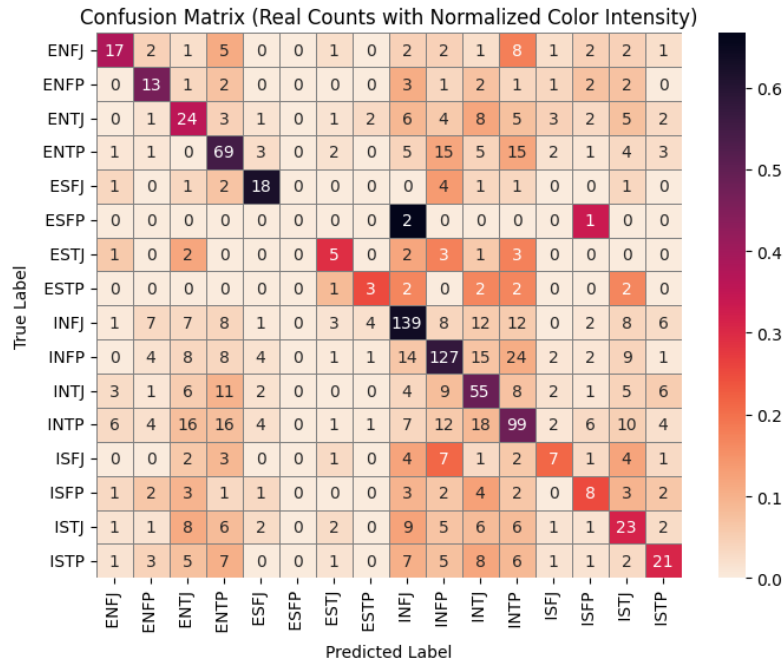


**Fig. 4** Confusion matrix for RoBERTa on MBTI personality classification.

## 6.3 MBTI Subtype Classification

Given RoBERTa's superior performance in MBTI classification, we further evaluate its ability to classify individual MBTI dimensions (E/I, S/N, T/F, J/P). Table 3 presents the F1-scores for each dichotomy, showing that the model achieves competitive performance across all four personality dimensions.

| MBTI Dichotomy | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Extraversion / Introversion (E/I) | 79% | 77% | 79% | 0.77 |
| Sensing / Intuition (S/N) | 81% | 78% | 81% | 0.78 |
| Thinking / Feeling (T/F) | 71% | 71% | 71% | 0.70 |
| Judging / Perceiving (J/P) | 71% | 71% | 71% | 0.71 |

**Table 3** Performance of RoBERTa for MBTI subtype classification.

## 6.4 Gender Classification

We also evaluated RoBERTa on the gender classification task. Table 4 summarizes the precision, recall, and F1-score for classifying gender from textual data. RoBERTa achieved an accuracy of **74.40%**, demonstrating balanced performance across both classes.

| Gender | Precision | Recall | F1-Score |
|---|---|---|---|
| Female | 73% | 77% | 0.75 |
| Male | 76% | 72% | 0.74 |
| Macro Avg. | 75% | 74% | 0.74 |
| Weighted Avg. | 75% | 74% | 0.74 |

**Table 4** Classification performance of RoBERTa on gender prediction.

The confusion matrix shown in Figure 5 further clarifies the classification behavior of RoBERTa. The model correctly identified 77% of females and 72% of males, with a relatively balanced distribution of misclassifications between both genders. This suggests RoBERTa effectively captures gender-specific linguistic patterns but indicates potential areas of overlap where communication styles between genders are not distinctly separable.

Overall, these results confirm RoBERTa's capability in distinguishing gender-based linguistic nuances within text data, although further improvement in distinguishing subtle differences remains necessary.

## 6.5 Analysis of Model Performance

Because of its optimal training schedule—which incorporates dynamic masking, bigger pretraining datasets, and the deletion of the next-sentence prediction objective—which
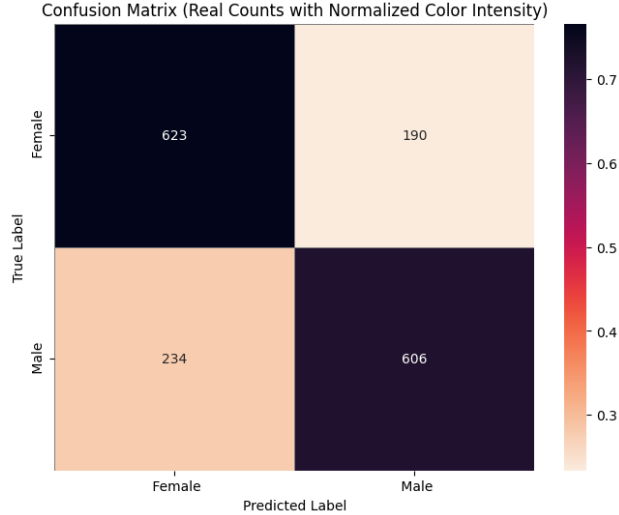
**Fig. 5** Confusion matrix for RoBERTa on gender classification.

so improves generalization-RoBERTa shows better performance across all classification tasks. Unlike GPT-2, which emphasizes generative tasks, RoBERTa models contextual relationships necessary for categorization by means of bidirectional encoding.

The subtype classification task demonstrates superior performance relative to the exact MBTI type classification. This corresponds with the hierarchical structure of MBTI traits, wherein categorizing broad personality dimensions (e.g., E/I) is fundamentally more straightforward than forecasting a 16-class outcome. Furthermore, gender classification demonstrates superior accuracy compared to MBTI classification, likely attributable to more pronounced linguistic differences linked to gender.

# 7 Conclusion

This study examines the effectiveness of Transformer-based language models, particularly RoBERTa, in classifying MBTI personality types and gender through conversational textual data sourced from Telegram forums. The findings demonstrate the considerable potential of utilizing linguistic cues from informal conversational contexts for implicit personality and demographic profiling.

Our findings indicate that RoBERTa successfully identifies personality-specific linguistic features, outperforming conventional methods in MBTI personality classification. The implementation of confidence thresholds emphasized the model's capacity to enhance prediction accuracy while sacrificing coverage, illustrating a significant trade-off relevant in precision-sensitive contexts.

In gender classification, RoBERTa demonstrated proficiency by effectively distinguishing gender-specific linguistic patterns, exhibiting balanced performance across classes. The analysis of the confusion matrix revealed distinct linguistic differences between male and female users, though some overlap in language patterns indicates that further refinement may enhance model differentiation.

Our analysis of MBTI feature distributions in conversational text highlights the substantial impact of personality-driven communication preferences on user interactions in digital platforms. The prevalence of introverted and intuitive personality types highlights their suitability for online text-based communication, whereas the lesser presence of sensing, judging, and extroverted types indicates a preference for face-to-face, structured interactions rather than digital engagement.

Although RoBERTa has shown effectiveness, challenges persist in implicit personality classification within informal conversational data. Challenges including noisy and unstructured language, lack of explicit labeling, and ethical concerns related to user privacy and data utilization remain prevalent. To address these challenges, advanced preprocessing techniques, adaptive tokenization strategies, and ethical frameworks are necessary to ensure the responsible use of implicit psychological modeling.

Future research should conduct in-depth analyses of linguistic features associated with specific personality and gender dimensions, create advanced fine-tuning strategies specifically designed for conversational data, and examine ethical implications in greater detail.

# Declarations

The authors declare that no ethical guidelines were violated, and no personally identifiable information was accessed or disclosed in conducting this research.

## Ethical Considerations and Compliance

All data utilized in this study were collected in full compliance with the Telegram platform's Terms of Service. Users' privacy and confidentiality were strictly maintained through anonymization procedures.

## Consent and Data Usage

All participants implicitly consented to the use of their anonymized messages for analytical purposes through their agreement with Telegram's Terms of Service. No personally identifiable information was disclosed or analyzed during this study.

## Conflict of Interest

The authors declare no conflicts of interest regarding the publication of this research.

## Ethical Approval

This research adhered strictly to ethical guidelines regarding data handling and privacy.

# References

[1] Myers, I.B., Myers, P.B.: Gifts Differing: Understanding Personality Type (1980)

[2] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is All You Need. In: Advances in Neural Information Processing Systems, vol. 30, pp. 5998–6008 (2017). https://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf

[3] Devlin, J., Chang, M.-W., Lee, K., Toutanova, K.: BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding (2018) arXiv:1810.04805 [cs.CL]

[4] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V.: RoBERTa: A Robustly Optimized BERT Pretraining Approach (2019) arXiv:1907.11692 [cs.CL]

[5] Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I.: Language Models are Unsupervised Multitask Learners (2019). OpenAI Blog

[6] Matz, S.C., Kosinski, M., Nave, G., Stillwell, D.J.: Psychological Targeting as an Effective Approach to Digital Mass Persuasion. Proceedings of the National Academy of Sciences **114**(48), 12714–12719 (2017) https://doi.org/10.1073/pnas.1710966114

[7] De Choudhury, M., Gamon, M., Counts, S., Horvitz, E.: Predicting Depression via Social Media. Proceedings of the Seventh International AAAI Conference on Weblogs and Social Media, 128–137 (2013)

[8] Hinds, J., Joinson, A.: Contextual Sensitivity and the 'Informed' User: Exploring the Impact of Privacy Controls on the Privacy Paradox. Journal of Broadcasting & Electronic Media **64**(4), 592–614 (2020) https://doi.org/10.1080/08838151.2020.1834297

[9] Ashraf, N., Naz, S.: Enhancing MBTI Personality Prediction from Text Data with Advanced Word Embedding Technique. VFAST Transactions on Software Engineering (2024)

[10] Li, B., Che, W.: Can Large Language Models Understand You Better? An MBTI Personality Detection Dataset Aligned with Population Traits. ArXiv (2024) arXiv:2401.12345 [cs.CL]

[11] Vásquez, R., Ochoa Luna, J.E.: Transformer-based Approaches for Personality Detection using the MBTI Model. In: 2021 XLVII Latin American Computing Conference (CLEI) (2021)

[12] Santos, V., Paraboni, I.: Myers-Briggs personality classification from social media text using pre-trained language models. J. Univers. Comput. Sci. (2022)

[13] Kadambi, P.: Exploring Personality and Online Social Engagement: An Investigation of MBTI Users on Twitter. ArXiv (2021) arXiv:2105.09347 [cs.CL]

[14] Bin Tareaf, R.: MBTI BERT: A Transformer-Based Machine Learning Approach Using MBTI Model For Textual Inputs. In: 2022 IEEE 24th International Conference on High Performance Computing & Communications (HPCC) and Associated Conferences (2022)

[15] Julianda, A.R., Maharani, W.: Personality Detection on Reddit Using Distil-BERT. Jurnal RESTI (Rekayasa Sistem dan Teknologi Informasi) (2023)

[16] Vonitsanos, G., Mylonas, P.: Decoding Gender on Social Networks: An In-depth Analysis of Language in Online Discussions Using Natural Language Processing and Machine Learning. In: 2023 IEEE International Conference on Big Data (BigData) (2023)

[17] Arya, S., Nishitha D'Souza, J.: Prediction of MBTI with textual data using different pre-trained transformer models. In: 2023 Fourth International Conference on Smart Technologies in Computing, Electrical and Electronics (ICSTCEE) (2023)

[18] Harrington, R., Loffredo, D.A.: Mbti personality type and other factors that relate to preference for online versus face-to-face instruction. The Internet and Higher Education **13**(1-2), 89–95 (2010)