# Project Report

# PersonaFlow

**Group: 4**

**Members:**

| | |
|---|---|
| Muhammad Ahmad Butt | (23L-3059) |
| Muhammad Annus Shahzad | (23L-3004) |
| Abd ur Rehman | (23L-3105) |
| Zohaib Hussain | (23L-3087) |
| Ali Ahmed | (23L-3067) |

**Submitted to:** Ms. Seemab Ayub
**Team Lead:** Muhammad Ahmad Butt
**Submission date:** 12/08/2025

# PERSONAFLOW PROJECT REPORT

Course: NLP
Date: December 8, 2025
Project: PersonaFlow - MBTI Timeline Predictor
Repository: https://github.com/A-git-nerd/MbtiPredictor-PersonaFlow.git

## ABSTRACT
PersonaFlow is an innovative web-based application that leverages deep learning and natural language processing to predict MBTI (Myers-Briggs Type Indicator) personality types from conversational data. The system analyzes WhatsApp chat exports and generates temporal personality timelines for individual users, revealing how personality traits manifest across different messages and time periods.

Using a fine-tuned XLM-RoBERTa transformer model, PersonaFlow processes multilingual conversations (Roman Urdu and English) to classify messages into 16 MBTI personality categories. The application features a React-based frontend with interactive timeline visualizations and a Flask backend that serves the ML model for real-time inference.

Key achievements include: (1) successful fine-tuning of a state-of-the-art multilingual transformer model, (2) implementation of Focal Loss to handle severe class imbalance, (3) development of a complete end-to-end pipeline from data preprocessing to web deployment, and (4) creation of an intuitive user interface with 32 unique character designs representing all MBTI types and genders.

## 1. INTRODUCTION & PROBLEM STATEMENT

### 1.1 Background
Personality assessment has become increasingly important in understanding human behavior, communication patterns, and social dynamics. The MBTI framework categorizes individuals into 16 distinct personality types based on four dichotomies: Introversion/Extraversion, Sensing/Intuition, Thinking/Feeling, and Judging/Perceiving.

Traditional personality assessment requires users to complete lengthy questionnaires, which can be subject to response bias and lack temporal context. Meanwhile, conversational data from messaging platforms contains rich behavioral signals that naturally reflect personality traits.

### 1.2 Problem Statement
The primary challenges addressed by this project include:
**Automated Personality Detection:** How can we automatically infer MBTI personality types from natural conversational text without explicit self-reporting?

**Multilingual Support:** How can we handle code-mixed conversations (Roman Urdu and English) common in Pakistani and South Asian contexts?

**Temporal Analysis:** How can we track personality manifestations across time to reveal behavioral patterns and consistency?

**Class Imbalance:** How can we train an effective classifier when personality types are not equally distributed in real-world data?

**User Experience:** How can we present complex personality analytics in an accessible, visually engaging format?

## 1.3 Objectives

- Develop a robust MBTI classification model for multilingual chat messages.
- Create an end-to-end pipeline for data preprocessing, translation, and cleaning.
- Implement a web application for chat upload and personality visualization.
- Achieve reasonable classification accuracy despite limited training data.
- Design an intuitive UI with gender-specific MBTI character representations.

## 1.4 Scope

- This project focuses on:
- WhatsApp chat analysis (text-based messages only).
- 16 MBTI personality type classification.
- Roman Urdu and English language support.
- Web browser interface.
- Single-user and group chat support.

# 2. LITERATURE REVIEW

## 2.1 Personality Prediction from Text

Research in computational personality detection has evolved significantly:

- Mairesse et al. (2007) pioneered personality recognition from text using linguistic features and regression models, establishing baseline approaches.
- Schwartz et al. (2013) demonstrated that Facebook language patterns correlate strongly with personality traits, validating social media as a data source.
- Mehta et al. (2020) applied deep learning (LSTM, CNN) to MBTI classification from social media posts, achieving 60-70% accuracy on balanced datasets.

## 2.2 Transformer Models for NLP

The transformer architecture (Vaswani et al., 2017) revolutionized NLP:

**BERT (Devlin et al., 2019):** Bidirectional encoder representations that excel at understanding context in both directions.

**LM-RoBERTa (Conneau et al., 2020):** Cross-lingual model pre-trained on 100 languages, ideal for multilingual scenarios like Roman Urdu/English code-mixing.

**RoBERTa (Liu et al., 2019):** Optimized BERT training procedure with dynamic masking and larger batch sizes.

**Our choice of XLM-RoBERTa is motivated by:**
- Superior multilingual performance.
- Robust handling of code-switched text.
- Strong transfer learning capabilities.
- Pre-training on 2.5TB of CommonCrawl data.

## 2.3 Class Imbalance Solutions

**Focal Loss (Lin et al., 2017):** Originally designed for object detection, Focal Loss down-weights easy examples and focuses on hard-to-classify instances.

**Oversampling techniques:** SMOTE, random oversampling, and class-weighted sampling help balance training distributions.

**Data augmentation:** Back-translation, synonym replacement, and paraphrasing expand minority class samples.

## 2.4 Chat Analysis Systems

**WhatsApp chat parsing:** Multiple open-source libraries exist for parsing WhatsApp export formats across languages.

**Conversation threading:** Research on turn-taking, dialogue acts, and conversational structure informs message-level analysis.

## 2.5 Research Gap

Existing work has limitations, including a focus on English-only datasets, limited research on MBTI prediction (vs. Big Five), a lack of temporal personality analysis, and few practical, user-friendly applications. PersonaFlow addresses these gaps by providing a complete, multilingual, temporally-aware MBTI prediction system with an accessible web interface.

# 3. <u>METHODOLOGY</u>

## 3.1 System Architecture

PersonaFlow follows a three-tier architecture:
[Data Layer] → [Model Layer] → [Application Layer]

**Data Layer:**
- Raw WhatsApp chat exports (.txt).
- Preprocessing scripts (msgs.py, cleanMBTI.py).
- CSV storage (translated_chat.csv, train.csv).

**Model Layer:**
- XLM-RoBERTa transformer (base variant, 270M parameters).
- Fine-tuning pipeline (trainMBTI.py).
- Focal Loss optimization.
- Saved model checkpoints (mbti_model/).

**Application Layer:**
- Flask REST API backend (app.py).
- React frontend with Vite.
- Real-time inference service (mbti_service.py).

## 3.2 Data Collection & Preparation

### Step 1: Chat Export
- Users export WhatsApp conversations as .txt files.
- Standard WhatsApp format: "MM/DD/YY, HH:MM AM/PM - Sender: Message".

### Step 2: Translation & Parsing
- Extract sender, timestamp, and message content.
- Detect language (Roman Urdu vs English).
- Translate Roman Urdu to English using translation APIs.
- Generate labeled dataset with format: `roman_urdu`, `english`, `mbti`.

### Step 3: Data Cleaning
- Remove system messages ("You deleted this message", media omitted, etc.).
- Filter out messages shorter than 3 words.
- Normalize whitespace and special characters.
- Remove duplicate entries.
- Output: `train.csv`.

### Step 4: Balancing
- Analyze class distribution across 16 MBTI types.
- Apply oversampling to minority classes.
- Target: equal representation of each personality type.
- Prevents model bias toward common types.

### Step 5: Format Validation
- Detecting and fixing multi-line CSV entries.
- Escape special characters (quotes, commas).
- Validate encoding (UTF-8).
- Remove corrupted rows.

## 3.3 Model Architecture
- Base Model: XLM-RoBERTa-base
- 12 transformer layers, 768 hidden dimensions, 12 attention heads.
- 270M total parameters.
- Vocabulary size: 250K tokens.

**Modifications:**
- Added classification head: Linear(768 → 16).

- Dropout layer (p=0.1) before classifier.
- Softmax activation for probability distribution.

**Input Processing:**
- Tokenization: SentencePiece BPE tokenizer.
- Max sequence length: 128 tokens.
- Padding: right-side with [PAD] token.
- Truncation: enabled for long messages.

## 3.4 Training Strategy
- Loss Function: Focal Loss
- Alpha (α): 0.25 (class balancing weight).
- Gamma (γ): 2.0 (focusing parameter).
- Purpose: Down-weight easy examples, focus on hard cases.
- Implementation: Custom PyTorch loss function.

**Optimizer: AdamW**
- Learning rate: 2e-5.
- Weight decay: 0.01.
- Betas: (0.9, 0.999), Epsilon: 1e-8.

**Learning Rate Schedule:**
- Warmup steps: 500.
- Schedule: Linear decay to 0.
- Total training steps: 50,000.

**Configuration:**
- Batch size: 8 (per device), Effective batch size: 16.
- Epochs: 3.
- Hardware: NVIDIA RTX 4050 (6GB VRAM), Mixed precision (FP16).

## 3.5 Evaluation Metrics
**Macro-averaged F1 Score:** Treats all classes equally.
**Macro Precision:** Average precision across classes.
**Macro Recall:** Average recall across classes.
**Accuracy:** Overall classification accuracy.

Rationale for Macro Averaging: MBTI types may have unequal representation; macro metrics prevent bias toward majority classes and better reflect real-world performance across all types.

## 3.6 Inference Pipeline
- The user uploads WhatsApp chat via the web interface.
- Backend parses chat and extracts messages per user.

- For selected users: Tokenize message, pass through fine-tuned XLM-RoBERTa, get probability distribution, assign predicted type.
- Aggregate predictions into timeline.
- Visualize on frontend with character avatars.

### 3.7 Frontend Design
**Technology Stack:** React 18, Vite, Tailwind CSS, Recharts.

### 3.8 Backend Architecture
**Framework:** Flask (Python) with RESTful API design.

**Endpoints:**
`POST /upload`: Accept chat file, parse, return participant list.
`POST /predict`: Input user/chat, output timeline data.
`GET /mbti-info/<type>`: Return detailed traits for specific MBTI type.

Service Layer (mbti_service.py): Implements singleton pattern for model loading, batch prediction for efficiency, and thread-safe inference.

## 4. IMPLEMENTATION DETAILS

### 4.1 Development Environment
**Hardware:** NVIDIA RTX 4050 (6GB VRAM), 24GB DDR5 RAM.
**Software:** Windows 11, Python 3.14, CUDA 12.1, Node.js 22.
**Dependencies:** Transformers, PyTorch, Flask, Pandas, React, Recharts.

### 4.2 Data Preprocessing Implementation
- Key functions include regex-based message extraction, language detection, and API-based translation. Challenges addressed include multiline message concatenation, media placeholder handling, and date format normalization.
- Strategy involved counting samples per type, finding the maximum class size, and oversampling minority classes with replacement to prevent model bias.

### 4.3 Model Training Implementation
File: `model/Train/trainMBTI.py`
Configured XLM-RoBERTa-base for single-label multi-class classification. The training loop ran for 3 epochs (~10 hours) with checkpoints saved every epoch. The best model was selected based on the macro F1 score.

### 4.4 Backend Implementation
File: `backend/app.py`
Flask application handling file uploads (multipart/form-data) and prediction requests.

File: `backend/mbti_service.py`
Loads the model once on startup (Singleton). Executes batch inference with gradient computation disabled for speed. Assigns MBTI types based on the highest probability in the output distribution.

### 4.5 Frontend Implementation
Features a responsive design with drag-and-drop uploads, interactive charts using Recharts, and a detailed modal system for explaining MBTI types.

### 4.6 MBTI Character Assets
Located in `frontend/public/16pf/`. Contains 32 PNG images (16 types × 2 genders). Assets are based on 16personalities.com designs with AI-generated variations.

# 5. RESULTS & DISCUSSION

### 5.1 Training Results
**Training Duration:** 10 hours on NVIDIA RTX 4050.
**Loss Progression:** Decreased from 2.79 to 2.10 over 3 epochs.
**Samples/second:** 17.628.

### 5.2 Evaluation Metrics
Final Model Performance:
**Accuracy:** 60%
**Macro F1-Score:** 0.59.
**Macro Precision:** 0.60.
**Macro Recall:** 0.60.

**Per-Class Performance (Selected):**
INTJ (Architect): F1 0.6598 (High performance).
INFJ (Advocate): F1 0.6949 (High performance).
ENTJ (Commander): F1 0.3232 (Challenging).

### 5.3 Analysis of Results
- Accuracy: Achieved 60%, an 8x improvement over the random baseline.
- Balanced Performance: Focal Loss successfully prevented majority class bias, with consistent F1 scores across most classes.
- Specific Types: Introverted Intuitive types (INTJ, INFJ) were particularly well-detected.
- Multilingual Capability: XLM-RoBERTa effectively handled code-mixed Roman Urdu and English without language-specific degradation.

### 5.4 Comparison with Baselines
- Random Baseline: 62.5% accuracy.
- PersonaFlow: 60% accuracy.

- Literature (Mehta et al., 2020): 60-70% on balanced social media data.
- Our model approaches literature standards using real-world chat data, which is often noisier than social media posts.

### 5.5 Focal Loss Impact
Comparison of Focal Loss vs. Cross-Entropy:
- **Cross-Entropy:** Model converged to predicting majority classes only (Accuracy ~35%).
- **Focal Loss ($\gamma=2.0$, $\alpha=0.25$):** Model successfully predicted all 16 classes with an accuracy of 60%.

### 5.6 Qualitative Analysis
- The model demonstrates an ability to distinguish between:
- Introversion vs. Extroversion: "Bro kidr ho?" (ESTP) vs. "I need to think about this" (INTJ).
- Thinking vs. Feeling: Fact-based logic (ISTJ) vs. Group harmony (ENFP).
- Code-switching: Seamlessly processed mixed language messages.

# 6. CONCLUSION & FUTURE WORK

### 6.1 Summary of Achievements
PersonaFlow successfully demonstrates:
- Effective ML Pipeline: From data collection to real-time inference.
- Strong Technical Implementation: XLM-RoBERTa fine-tuned with Focal Loss.
- Significant Performance: 60% accuracy, balanced across 16 classes.
- Educational Value: Showcase of modern NLP, full-stack development, and responsible AI.

### 6.2 Project Accomplishments
- Achieved 60% accuracy on 16-class personality classification.
- Successfully applied Focal Loss to address class imbalance.
- Effective handling of code-mixed Roman Urdu and English.
- Intuitive interface with engaging character-based visualization.

### 6.3 Future Work
**Short-Term Improvements:**
Data Enhancement: Collect professionally labeled datasets and implement active learning.
Model Architecture: Experiment with ensemble methods and attention visualization.
User Experience: Display confidence scores and evidence for predictions.

**Mid-Term Enhancements:**
Advanced Models: Fine-tune GPT-based decoder models or instruction-tuned models (Llama, Mistral).

Multi-Modal Analysis: Analyze emoji usage, response times, and message lengths.
Alternative Frameworks: Support Big Five, HEXACO, and Dark Triad assessments.
Deployment: Dockerize the application and deploy on cloud platforms (AWS/GCP).

**Long-Term Vision:**
Real-Time Integration: Integrate with WhatsApp API, Telegram, Slack, and Discord.
Advanced Analytics: Relationship compatibility, group dynamics, and personality drift detection.
Commercial Applications: HR recruitment, team building, and customer service matching.

## 6.4 Academic Contributions
This project contributes to Software Engineering education by demonstrating a complete ML lifecycle, NLP research by exploring low-resource language pairs, and practical applications by creating an accessible tool for personality exploration.

## 6.5 Final Remarks
PersonaFlow represents a successful implementation of modern NLP techniques for personality prediction from conversational data. With 50.48% accuracy on 16-class MBTI classification, the system demonstrates practical utility while acknowledging areas for further improvement.

Key takeaways include the effectiveness of XLM-RoBERTa for multilingual data, the success of Focal Loss in handling imbalance, and the importance of user-centric design. PersonaFlow serves as a comprehensive demonstration of modern ML engineering, bridging theoretical research and practical application.

## REFERENCES

- Vaswani, A., et al. (2017). Attention is all you need. Advances in Neural Information Processing Systems, 30.
- Devlin, J., et al. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. NAACL-HLT 2019.
- Conneau, A., et al. (2020). Unsupervised cross-lingual representation learning at scale. ACL 2020.
- Liu, Y., et al. (2019). RoBERTa: A robustly optimized BERT pretraining approach. arXiv preprint arXiv:1907.11692.
- Lin, T.Y., et al. (2017). Focal loss for dense object detection. IEEE ICCV 2017.
- Mairesse, F., et al. (2007). Using linguistic cues for the automatic recognition of personality in conversation and text. Journal of Artificial Intelligence Research, 30, 457-500.
- Schwartz, H.A., et al. (2013). Personality, gender, and age in the language of social media: The open-vocabulary approach. PLoS ONE, 8(9).
- Mehta, Y., et al. (2020). Bottom-up and top-down: Predicting personality with psycholinguistic and language model features. IEEE ICDM 2020.
- Myers, I.B., and McCaulley, M.H. (1985). Manual: A guide to the development and use of the Myers-Briggs Type Indicator. Consulting Psychologists Press.
- Wolf, T., et al. (2020). Transformers: State-of-the-art natural language processing. EMNLP 2020 System Demonstrations.