

# NLP Graphic Tool: Research Assignment Option 2

## Natural Language Processing Course

Tesfay Hagos Weldegebriel (VR544201)  
tesfayhagos.weldegebriel@studenti.univr.it

February 27, 2026

## Course Information

- **University:** University of Verona
  - **Professor:** Prof. Matteo Cristani
  - **Subject:** Natural Language Processing
- 

## 1 Introduction

This report documents the implementation of a **Graphic Tool for Document Processing**, developed as part of the Research-Oriented Assignment (Option 2) for the Natural Language Processing course. The tool is designed to provide a modern, interactive interface for common text analysis tasks, with specialized support for both **English** and **Tigrinya**.

**Live Application:** The tool is deployed and accessible at:  
<https://weldegebriel-tesfayhagos-nlp.streamlit.app/>

## 2 System Architecture

The application is built using a modular Python architecture:

- **Core Pipeline (nlp\_pipeline.py):** Handles text cleaning, language detection, tokenization, lemmatization (English/Tigrinya), and relevance calculations.
- **Graphic UI (app.py):** An interactive web interface built with **Streamlit**, featuring real-time analysis, visualizations with **Plotly**, and data export capabilities.

## 3 User Interface

The application features a modern, responsive interface designed for ease of use. Below are screenshots illustrating the main dashboard and analysis results.

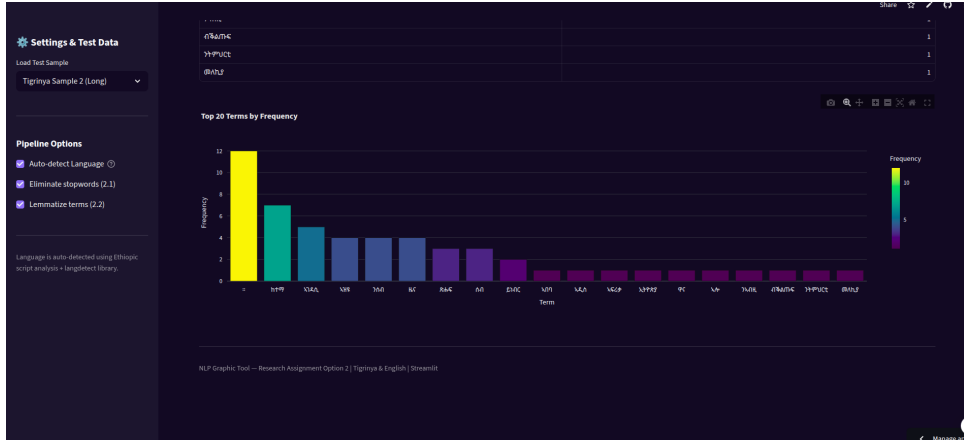


Figure 1: Analysis Overview: Dashboard showing metrics and token flow.



Figure 2: Interactive Insights: Visualizations of term frequencies and relevance.

## 4 Requirement Implementation

The project successfully addresses all five requirements specified in the assignment:

### 4.1 2.1 Eliminate Stopwords

The tool removes common stopwords that carry little semantic weight.

- **English:** Uses the standard NLTK stopwords corpus.
- **Tigrinya:** Integrates with the `tigrinya-nlp` library using a minimal stopwords configuration.

### 4.2 2.2 Lemmatization

- **English:** Implements the `WordNetLemmatizer` with POS (Part-of-Speech) tagging for high accuracy.
- **Tigrinya:** Implements a custom **rule-based stemmer** developed with verification from a native speaker. It handles possessive, plural, and object pronoun suffixes (e.g., suffix stripping (e.g., stripping plural and pronoun suffixes)).

## 4.3 2.3 Frequency Computation

The pipeline generates a frequency table of all lemmas/stems after filtering. Results are displayed in interactive tables and visualized via bar charts and pie charts.

## 4.4 2.4 Distance from Strategic Points

Distances are measured relative to the total number of tokens:

1. **Distance from Start:**  $\frac{\text{First occurrence position}}{\text{Total tokens}}$
2. **Distance from End:**  $\frac{\text{Total tokens} - \text{Last occurrence position}}{\text{Total tokens}}$

## 4.5 2.5 Compound Relevance Indices

A custom relevance index is calculated following the 50% frequency + 50% earliness formula:

$$\text{Relevance} = 0.5 \times \text{Frequency Score} + 0.5 \times \text{Earliness Score}$$

Where:

- *Frequency Score* =  $\frac{\text{Term Frequency}}{\text{Max Frequency}}$
- *Earliness Score* =  $1 - \frac{\text{Average Position}}{\text{Total Tokens}}$

## 5 Key Features

- **Language Detection:** Automatic detection between English and Tigrinya.
- **Interactive Visualizations:** Scatter plots (Distance vs. Pos), Bar charts (Relevance), and Pie charts (Distribution).
- **Data Export:** Analysis results can be downloaded in **CSV** or **JSON** formats.
- **Modern UI:** Features a dark-themed glassmorphic design for enhanced user experience.

## 6 Conclusion

The NLP Graphic Tool provides a robust framework for document analysis, fulfilling academic requirements while pushing beyond basic processing by supporting a low-resource language like Tigrinya with specialized linguistic rules.