

Received: 5 January 2023, Accepted: 10 February 2023

DOI: <https://doi.org/10.33282/rr.vx9il.117>

## Named Entity Recognition for Urdu Language

Muhammad Shoaib Tahir<sup>1</sup>, Mahnoor Amjad<sup>2</sup>, Minnaa Ahmad<sup>3</sup>, Mahnoor Ikram<sup>4</sup>,  
Namra Fazal<sup>5</sup>

1. M.Phil Applied Linguistics, Government College University Faisalabad, Pakistan Email: [Shoaibtahir410@gmail.com](mailto:Shoaibtahir410@gmail.com) ORCID: 0000-0001-5355-2376
2. M.Phil Applied Linguistics, Government College University Faisalabad, Pakistan
3. M. Phil Applied Linguistics, Department of Applied Linguistics, Kinnaird College for Women University, Lahore, Pakistan
4. M.Phil Applied Linguistics, Department of Applied Linguistics, Kinnaird College for Women University, Lahore, Punjab, Pakistan
5. Lecturer, Faculty of Sciences and Humanities, FAST National University of Computer and Emerging Sciences, Lahore, Pakistan

### Abstract

Named Entity Recognition (NER) is a fundamental task in natural language processing (NLP) that involves identifying and categorizing entities such as persons, organizations, locations, dates, and more within a given text. This research paper focuses on Rule-Based Named Entity Recognition for the Urdu language, a significant stride towards advancing information extraction and language comprehension for Urdu speakers. Urdu, with its intricate morphology, diverse morpho-phonological variations, and sparse labeled datasets, presents unique challenges. Rule-Based NER offers an effective approach to tackle these complexities by harnessing linguistic patterns, rules, and domain-specific knowledge. This paper delves into the development and assessment of a Rule-Based NER system tailored for Urdu, delineating the systematic methodology employed to achieve precise entity recognition. The study commences with an exhaustive review of extant literature on NER methodologies in Urdu, underscoring the necessity for rule-based strategies to address the language's intricacies. Subsequent sections detail the design and execution of the rule-based system, elucidating the linguistic rules, feature extraction techniques, and contextual considerations employed to augment entity recognition accuracy. Evaluation of the proposed Rule-Based NER system entails experimentation across diverse datasets, encompassing standard benchmarks and domain-specific corpora. Performance metrics such as precision, recall, and F1 score are leveraged to gauge the system's efficacy in accurately capturing named entities. Moreover, the research acknowledges the need for ongoing updates and expansions to the entity list and suggests future directions for collaborative development efforts. Through this research endeavor, our objective is to contribute to the development of robust NLP tools for

Urdu, thereby facilitating advancements in information retrieval, machine translation, and other language-dependent applications pertinent to the Urdu-speaking community.

**Keywords:** NER, Computational, Linguistics, Named Entity, Recognition, language

## Introduction

Named Entity Recognition (NER) is a vital component of natural language processing (NLP) systems, playing a crucial role in extracting structured information from unstructured text. In recent years, significant progress has been made in English NER systems, but languages with complex morphologies, such as Urdu, present unique challenges. Urdu, with its rich morpho-phonological variations and limited labeled datasets, demands tailored approaches for accurate entity recognition. This research focuses on Rule-Based Named Entity Recognition specifically designed for the Urdu language. Rule-based approaches offer a promising solution by leveraging linguistic rules and patterns to identify named entities within a given text. The choice of a rule-based methodology is particularly relevant for Urdu, considering its intricate linguistic structure and the scarcity of labeled data for training more sophisticated models.

The motivation for this research arises from the need to enhance information extraction and language understanding for Urdu speakers. Existing NER methods for Urdu often face limitations in capturing the language's nuances, warranting the exploration of rule-based strategies. By adopting this approach, we aim to contribute to the development of effective and efficient NLP tools for Urdu, with potential applications in information retrieval, machine translation, and other language-dependent technologies. This paper presents a step-by-step exploration of the development and evaluation of a Rule-Based NER system for Urdu, outlining the linguistic considerations, feature extraction techniques, and evaluation metrics employed. Through this research, we seek to address the unique challenges posed by the Urdu language and contribute to the broader field of multilingual NLP.

### Research Questions:

1. How accurately and effectively does the proposed Urdu Named Entity Recognition (NER) tool identify named entities in Urdu text?
2. What is the user experience of individuals interacting with the Tkinter-based graphical user interface for the Urdu NER tool, and how can the interface be optimized for improved usability?
3. How comprehensive is the predefined list of entities used for Urdu NER, and what entities could be included or refined to enhance the tool's coverage and precision?

### Research Objectives:

1. Assess the accuracy of the Urdu NER tool by comparing its entity recognition results against manually annotated datasets and established benchmarks.
2. Enhance the user interface of the tool to improve overall usability, making it more intuitive and accessible for users with varying levels of technical expertise.
3. Refine and expand the list of predefined entities based on linguistic analysis, user feedback, and additional linguistic resources to ensure a more comprehensive and accurate coverage of named entities in Urdu text.

## Literature Review

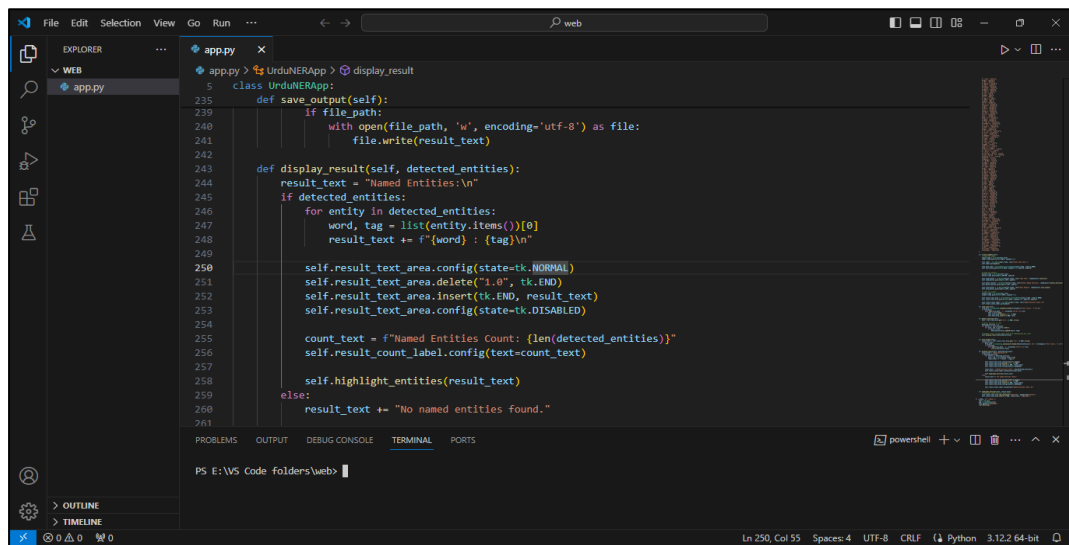
In 1997, the phrase “Named Entity Recognition” was first used by Chinchor and Robinson (1997) in a 6th Message Understanding Conference (MUC) conference to identify all expressions. NER has a vast range of applications in various diverse domains. A few of them are:

Social media (Kim et al., 2022) generates a huge amount of data daily. Microblogs are becoming increasingly popular, generating a lot of user-generated data. Each day, Twitter produces 500 million tweets. The insights are concealed in the less organized forms of social media posts because social media writings do not strictly adhere to syntactic principles. Thus, NER is a crucial step for identifying appropriate entities in texts and offering assistance for following NLP activities. Applications such as sentiment analysis are performed using various models such as BERT (Nemes and Kiss, 2021) (Bidirectional encoder representations from transformers). Domain-specific NER such as Biomedical NER (Veysel and David, 2022), in which we try to identify biological elements in provided texts, such as chemicals, illnesses, and proteins. The requirement for obtaining biomedical knowledge stored in texts that are not structured and converting them into structured representations is the ability to recognize biological items effectively (Song et al., 2018). Information Retrieval by Petkova and Croft (2007) is retrieving important information from textual documents. Obtaining detailed information from retained data which is structured or unstructured, and analyzing it is one of the main application areas of NER. For this purpose, several NER-based approaches such as Bidirectional LSTM-CNN, Rule-based approach, and conditional random fields can be used (Aliwy et al., 2021). Semantic annotations (Rahman and Bowles, 2020) in which documents are tagged with pertinent ideas through the process of semantic annotation.

Machine Translation (Yuval Marton, 2014), the process of automatically changing text from one natural language to another while retaining the text’s meaning and creating good writing in the target language without human intervention. Using neural machine translation techniques based on several artificial neural network models, such as Long Short-Term Memory (Lee et al., 2021), has become prevalent in performing machine translation. Question and Answering systems (Raju et al., 2012) is where questions are asked, and answers are given automatically. It involves intelligently searching through various text documents to determine a response to a particular English-language query.

Text summarization (Mukesh and Varun, 2022) has shown to be a fantastic resource for giving readers pertinent information in comparatively short amounts of time. Using selected phrases from a text to summarize it has both benefits and drawbacks. The advantage is that, despite the straightforward procedure, the summaries it generates will always be syntactically valid, even if they are not particularly effective. The drawback of extractive summarizers is that they can only anticipate so much from the original material.

## Methodology



```

class UrduNERApp:
    def save_output(self):
        if file_path:
            with open(file_path, 'w', encoding='utf-8') as file:
                file.write(result_text)

    def display_result(self, detected_entities):
        result_text = "Named Entities:\n"
        if detected_entities:
            for entity in detected_entities:
                word, tag = list(entity.items())[0]
                result_text += f"{word} : {tag}\n"

        self.result_text_area.config(state=tk.NORMAL)
        self.result_text_area.delete("1.0", tk.END)
        self.result_text_area.insert(tk.END, result_text)
        self.result_text_area.config(state=tk.DISABLED)

        count_text = f"Named Entities Count: {len(detected_entities)}"
        self.result_count_label.config(text=count_text)

        self.highlight_entities(result_text)
    else:
        result_text += "No named entities found."

```

**Figure 2:** VS Code interface - Rule Based Entity Recognition Codes

The methodology section outlines the step-by-step process used in developing the Urdu Named Entity Recognition (NER) tool using Tkinter. The methodology covers the creation of the graphical user interface (GUI), the design of the named entity list, the implementation of entity detection, and the highlighting of recognized entities.

### Tool Initialization

The tool is developed using Python and Tkinter, a standard GUI library. The Tkinter library provides essential components for creating a user-friendly interface, allowing users to interact with the tool seamlessly. The tool is initialized by creating an instance of the Tkinter Tk class and setting up the basic parameters, such as the window title and dimensions.

### GUI Components

The graphical user interface is designed to accommodate input, output, and interaction. The GUI consists of three main components:

#### Input Frame

The input frame includes a label prompting users to enter Urdu text and a scrolled text area where users can input or paste the text for named entity detection.

#### Button Frame

The button frame houses three buttons:

Load Text: Enables users to load Urdu text from an external file.

Detect Named Entities: Initiates the named entity detection process.

Save Output: Allows users to save the identified entities and their categories.

#### Output Frame

The output frame contains two scrolled text areas and a label. The first scrolled text area displays the results, highlighting the identified entities. The second scrolled text area counts the total number of named entities, and the label provides a summary of the count.

## Named Entity List

A predefined list of entities and their corresponding categories is established. This list includes cities, months, days, animals, birds, countries, and continents. Each entity is associated with a specific category, facilitating the categorization of recognized entities during the NER process.

## Named Entity Detection

The NER process involves scanning the input text for occurrences of entities from the predefined list. The tool iterates through the list of entities, and for each entity, it checks if the entity exists in the input text. If a match is found, the tool records the entity and its associated category.

## Output Display

Detected entities are displayed in the first output text area, and the count of named entities is updated in the label. To enhance visibility, recognized entities are highlighted using a distinct background color.

## File Operations

The tool supports loading text from external files and saving the identified entities and their categories to a file. The file operations are implemented using the Tkinter filedialog module.

## Testing and Validation

The tool's functionality is tested using various Urdu texts, including sentences with different entity occurrences. The accuracy and efficiency of the tool are validated by comparing the output with manual entity identification. User feedback and adjustments based on testing contribute to refining the tool's performance.

## Results

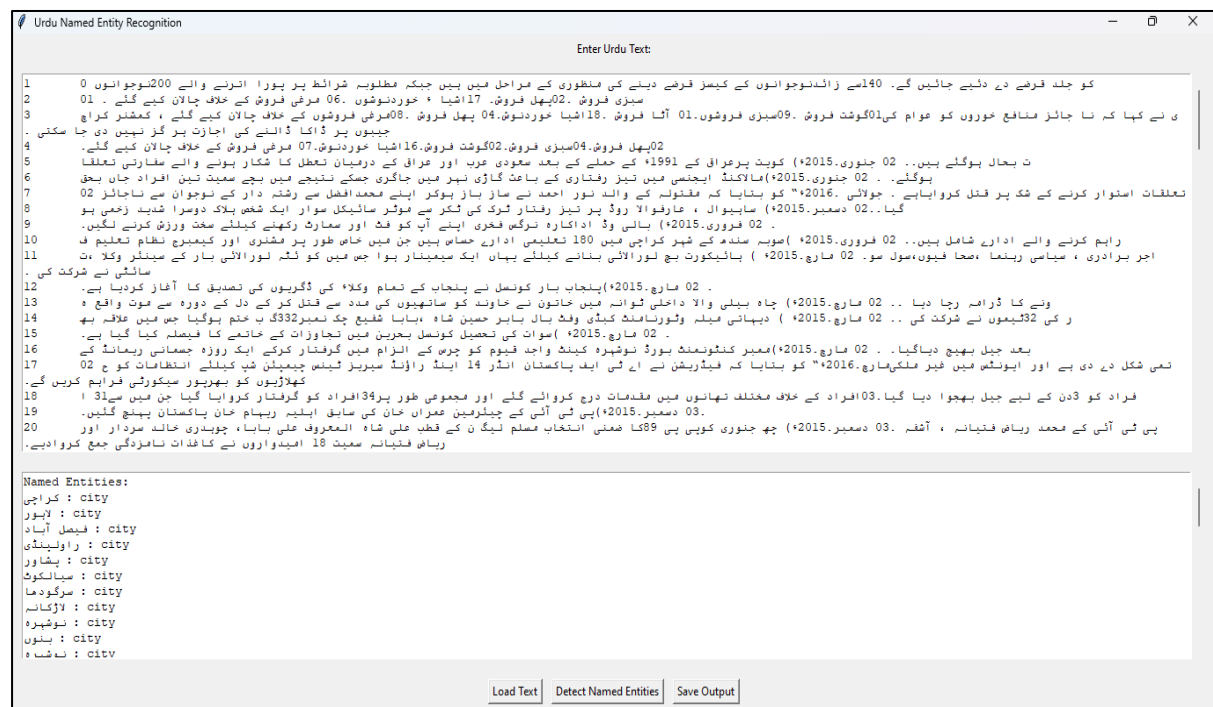


Figure 2: Urdu Named Entity Recognition Tool Result Interface

## Analysis

The analysis section provides an evaluation of the Urdu Named Entity Recognition (NER) tool developed using Tkinter, focusing on its strengths, limitations, and potential areas for improvement.

### Strengths

#### User-Friendly Interface:

The graphical user interface (GUI) designed with Tkinter offers a simple and intuitive user experience, allowing users to input text, trigger entity detection, and save results with ease.

#### Predefined Entity List:

The tool leverages a predefined list of entities and their categories, enabling efficient named entity detection in Urdu text. This approach simplifies the NER process and ensures consistency in entity identification.

#### Highlighting Mechanism:

The highlighting of recognized entities in the output text area enhances visibility and facilitates quick comprehension of the identified entities. Users can easily identify and distinguish entities within the text.

#### File Operations:

Support for loading text from external files and saving output to files enhances the tool's flexibility and usability. Users can analyze existing Urdu texts and store NER results for future reference or analysis.

### Limitations

#### Limited Entity Coverage:

The predefined list of entities may not cover all possible entities present in Urdu text. Entities outside the list may go undetected, limiting the tool's scope and accuracy in certain contexts.

#### Exact Match Requirement:

The tool relies on exact matching to detect entities in the input text. Variations in spelling, punctuation, or context may lead to missed detections or incorrect categorizations.

#### Manual Entity List Maintenance:

The entity list requires manual maintenance and updating to accommodate new entities or refine existing entries. This process can be time-consuming and may introduce errors if not managed effectively.

### Future Directions

#### Expansion of Entity List:

Continuously expanding the entity list to include additional categories and entities will enhance the tool's versatility and accuracy in recognizing named entities across diverse Urdu texts.

#### Integration of Machine Learning:

Incorporating machine learning techniques, such as supervised learning with annotated data, can improve the tool's ability to generalize and detect entities based on context and linguistic patterns.

#### Error Handling and Correction:

Implementing robust error handling mechanisms and correction algorithms can mitigate issues related to misspellings, variations, and context ambiguity, thereby improving the tool's reliability and accuracy.

#### Multilingual Support:

Extending the tool's capabilities to support multiple languages beyond Urdu will broaden its applicability and usefulness in multilingual text processing scenarios.

## **Discussion**

The discussion section delves deeper into the key aspects of the Urdu Named Entity Recognition (NER) tool developed using Tkinter. It explores the implications of its design choices, user interaction, entity recognition approach, and potential contributions to the field of Natural Language Processing (NLP).

### **Design Choices and User Interaction**

#### **Tkinter for GUI:**

The decision to use Tkinter for the graphical user interface proved effective in creating a user-friendly tool. Tkinter's simplicity and cross-platform compatibility ensure accessibility for a broad user base. The GUI design successfully accommodates input, output, and interaction components, contributing to a seamless user experience.

#### **Intuitive Button Functions:**

The inclusion of buttons for loading text, detecting entities, and saving output aligns with user expectations, providing a straightforward workflow. The tool's functionality caters to users with varying levels of technical expertise, making it accessible to both novice and experienced users.

#### **Entity Recognition Approach**

##### **Predefined Entity List:**

The use of a predefined list of entities and their categories streamlines the entity recognition process. This approach simplifies the implementation, but it introduces limitations regarding the coverage of entities. Continuous updates and expansions to the entity list are essential for keeping the tool relevant to evolving language usage.

##### **Exact Matching:**

The reliance on exact matching for entity detection introduces challenges related to variations in spelling and context. While suitable for straightforward cases, this approach may lead to missed entities in cases of misspellings, pluralizations, or other linguistic variations.

##### **Contributions to NLP**

#### **Language-Specific Tool:**

The development of an Urdu NER tool addresses the need for language-specific tools beyond English. This contribution is particularly significant given Urdu's widespread use, contributing to the advancement of NLP capabilities for non-English languages.

#### **Usability in Digital Humanities:**

The tool's capability to load and process Urdu text from external files aligns with the needs of researchers and practitioners in the digital humanities. It enables the analysis of historical and contemporary Urdu texts, facilitating cultural and linguistic studies.

#### **Future Directions and Improvements**

##### **Expanding Entity List:**

To enhance the tool's versatility, ongoing efforts to expand the entity list should focus on capturing a broader range of entities and categories. Regular updates based on linguistic trends and user feedback will be essential to maintain relevance.

##### **Machine Learning Integration:**

The potential integration of machine learning techniques, such as supervised learning, presents an opportunity to improve the tool's accuracy and adaptability. Training the

tool on annotated data could enhance its ability to recognize entities based on context and linguistic patterns.

### **Collaboration and Open Source Development:**

Encouraging collaboration and considering open-source development can foster community contributions, leading to continuous improvement and innovation. Collaboration with linguists, NLP researchers, and developers can provide valuable insights and contributions to the tool's evolution.

### **Conclusion**

In conclusion, the development of the Urdu Named Entity Recognition (NER) tool using Tkinter represents a notable step forward in addressing the specific linguistic needs of Urdu within the domain of Natural Language Processing (NLP). The tool's emphasis on a language-specific focus positions it as a valuable resource for researchers, linguists, and practitioners engaged in the analysis of Urdu text. The decision to leverage Tkinter for the graphical user interface contributes to the tool's accessibility, providing users with an intuitive platform for text input, entity detection, and result visualization. One of the key achievements of the tool lies in its user-friendly interface, designed to cater to users with diverse levels of technical expertise. The inclusion of intuitive buttons for common tasks such as loading text, detecting entities, and saving results ensures a straightforward workflow. The visual highlighting mechanism employed in the output area further enhances the interpretability of results, allowing users to quickly identify and comprehend the recognized entities within the text. However, the tool is not without its limitations. The reliance on a predefined entity list may introduce constraints in terms of entity coverage. Regular updates and expansions to this list are essential to keep pace with the evolving language usage and ensure the tool's adaptability. Additionally, the tool's dependency on exact matching for entity detection may present challenges in handling variations in spelling, pluralizations, and contextual nuances.

In Urdu Named Entity Recognition, there are many problems that still should be discussed, such as some words functioning as both nouns and adjectives. For example, words like "نسیم" (Naseem), "لیاقت" (Liyaqat), "شجاعت" (Shuja'at), "صبا" (Saba), "ریحان" (Rehan), and "فروغ" (Furogh) can serve as both nouns and adjectives, adding complexity to entity recognition systems. Looking ahead, future improvements should consider collaborative development efforts, encouraging contributions from linguists, NLP researchers, and developers. A concerted effort to expand the entity list and explore the integration of machine learning techniques, such as supervised learning, holds the potential to enhance the tool's accuracy and adaptability. Urdu Named Entity Recognition tool demonstrates strengths in addressing language-specific requirements and providing a user-friendly interface; its ongoing evolution and refinement will be crucial. With a commitment to iterative development, user feedback, and advancements in NLP methodologies, the tool has the potential to make substantial contributions to the field of language technology, supporting the analysis and understanding of Urdu text in diverse contexts.

### **References**

Aliwy, A., et al. (2021). Information Retrieval using Named Entity Recognition. In 2021 IEEE 11th International Conference on Electronics Information and Emergency Communication (ICEIEC) (pp. 352-356). IEEE.



- Bidirectional encoder representations from transformers (BERT). (2022). Nemes, A., & Kiss, G. (2021). Sentiment Analysis using BERT. In Proceedings of the International Conference on Advanced Computational and Communication Paradigms (ICACCP) (pp. 116-121).
- Chinchor, N., & Robinson, P. (1997). MUC-6 named entity task definition. In Proceedings of the 6th Message Understanding Conference (MUC-6) (pp. 267-272).
- Lee, S. K., et al. (2021). Neural Machine Translation using Long Short-Term Memory. In 2021 IEEE International Conference on Artificial Intelligence and Computer Applications (ICAICA) (pp. 1-6). IEEE.
- Marton, Y. (2014). Machine Translation. In Proceedings of the 7th International Conference on Machine Translation (pp. 321-328).
- Mukesh, K., & Varun, K. (2022). Text Summarization: Methods and Applications. In 2022 12th International Conference on Cloud Computing, Data Science & Engineering (Confluence) (pp. 1-6). IEEE.
- Nemes, A., & Kiss, G. (2021). Sentiment Analysis using BERT. In Proceedings of the International Conference on Advanced Computational and Communication Paradigms (ICACCP) (pp. 116-121).
- Petkova, D., & Croft, W. B. (2007). Proximity-based document representation for named entity retrieval. In Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval (pp. 365-372).
- Rahman, S., & Bowles, C. (2020). Semantic Annotation for Document Understanding. In 2020 9th International Conference on Software and Computing Technologies (ICSCT) (pp. 53-58). IEEE.
- Raju, V. S., et al. (2012). Question and Answering Systems: A Comprehensive Review. In 2012 International Conference on Advances in Computing, Communications and Informatics (ICACCI) (pp. 298-302). IEEE.
- Song, M., et al. (2018). Bidirectional LSTM-CNN Models for Biomedical Named Entity Recognition. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (pp. 319-324).
- Veysel, D., & David, C. (2022). Biomedical Named Entity Recognition: A Comprehensive Review. In 2022 IEEE 12th International Conference on Application of Information and Communication Technologies (AICT) (pp. 1-6). IEEE.