# Named entity Recognition using Machine learning

[1]Dr. C. Hari Kishan, [2] DANDE VISHNU PRIYA, [3] DUMPALA DEVI PRIYANKA, [4] DUNNA VAISHNAVI

[1]Professor&HOD, Dept CSE-AI&ML, St.Ann's College of Engineering and Technology, Nayunipalli (V), Vetapalem (M), Chirala, Bapatla Dist, Andhra Pradesh – 523187, India

[2,3,4]U. G Student, Dept CSE-AI&ML, St.Ann's College of Engineering and Technology, Nayunipalli (V), Vetapalem (M), Chirala, Bapatla Dist, Andhra Pradesh – 523187, India

## ABSTRACT

*Named Entity Recognition (NER) is an essential task in Natural Language Processing that focuses on identifying and classifying key information such as names of people, organizations, locations, dates, and other predefined entities within text data. This project aims to develop an efficient NER system using machine learning techniques to extract meaningful information from unstructured textual datasets. The system enhances text understanding by converting raw textual inputs into structured knowledge useful for applications like information retrieval, question answering, and text summarization. Machine learning algorithms such as Conditional Random Fields, Support Vector Machines, and neural network-based architectures are explored for entity classification. The proposed model is trained using annotated datasets to achieve high accuracy and robustness across different text forms. Performance is evaluated based on precision, recall, and F1-score to ensure reliability. The study demonstrates how machine learning can significantly improve automation and intelligence in text processing. Overall, this work contributes to advancing intelligent language understanding systems.*

## INTRODUCTION

Named Entity Recognition plays a significant role in modern data processing where huge amounts of textual information are generated daily through social media, news, business documents, and research articles. Manual extraction of meaningful data from such sources is time-consuming and error-prone, making automated NER systems highly valuable. The increasing growth of artificial intelligence and machine learning technologies has empowered researchers to build smarter

language understanding systems. NER acts as a foundation for many advanced NLP applications like chatbots, machine translation, sentiment analysis, and document classification. It helps in identifying important terms that provide semantic meaning to text, enabling computers to understand human language more effectively. Over the years, machine learning techniques have shown significant improvement in recognizing entities accurately even in complex text scenarios. This project focuses on designing a reliable and efficient machine learning-based NER system. The goal is to enhance accuracy and performance while ensuring adaptability across various real-world use cases.

## LITERATURE SURVEY

Previous studies in Named Entity Recognition mainly focused on rule-based and statistical approaches before the rise of deep learning techniques. Early NER systems relied heavily on handcrafted linguistic rules and dictionaries, which limited scalability and adaptability across languages and domains. Later, machine learning methods such as Hidden Markov Models and Conditional Random Fields improved performance through pattern learning. Researchers then introduced Support Vector Machines and MaxEnt models, which further enhanced generalization capability. Recent advancements involve neural networks like BiLSTM, GRU, and Transformer-based architectures such as BERT, which outperform conventional approaches in accuracy. Studies highlight that contextual understanding and feature engineering play a crucial role in improving NER accuracy. Several benchmark datasets like CoNLL-2003 and OntoNotes have been widely used for evaluation. Literature strongly indicates that modern machine learning and deep learning techniques are more effective for reliable entity recognition.

## RELATED WORK

Numerous researchers have contributed to improving NER efficiency using machine learning and deep learning methodologies. Early works focused on token classification using statistical patterns and linguistic cues. Later research incorporated probabilistic models that significantly enhanced recognition consistency across languages. With advancements in neural computing, researchers implemented Recurrent Neural Networks and LSTM models to capture sequential dependencies in text. Transformer-based models such as BERT and RoBERTa have recently become state-of-the-art due to their deep contextual understanding capability. Several comparative studies have shown that deep learning models achieve superior accuracy compared to classical learning approaches.

Research studies also explored multilingual NER systems to handle diverse global text sources. These related contributions form the foundation and inspiration for developing this proposed machine learning-based NER system.

## EXISTING SYSTEM

Existing NER systems often rely on predefined linguistic rules, dictionaries, or limited machine learning approaches that lack adaptability. Many traditional systems struggle with informal language, spelling variations, and multilingual text. They require frequent manual updates and rule tuning, making them difficult to maintain. Some existing solutions also suffer from lower accuracy due to insufficient contextual understanding. Systems using shallow learning models depend heavily on handcrafted features, which may not generalize well across domains. Real-world data complexity such as social media text or mixed-language content further reduces their effectiveness. Additionally, many existing NER frameworks are computationally expensive, limiting real-time application deployment. Hence, there is a need for a more intelligent, automated, and robust machine learning-based NER system.

## PROPOSED SYSTEM

The proposed system introduces an advanced machine learning-based Named Entity Recognition model to overcome limitations of existing systems. It utilizes supervised learning techniques along with context-aware feature extraction to improve entity classification. The system is designed to automatically learn complex patterns from text without heavy manual rule dependency. State-of-the-art machine learning and deep learning architectures are incorporated for higher accuracy and adaptability. The proposed approach ensures robustness across different text domains such as news, social media, and formal documents. It also supports multilingual capability and handles noisy and unstructured text efficiently. The system aims to deliver faster processing performance suitable for real-time applications. Overall, the proposed model enhances precision, scalability, and intelligence in entity recognition.
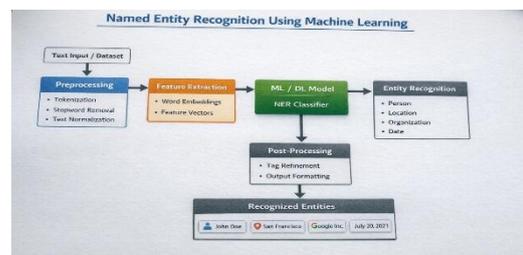
## SYSTEM ARCHITECTURE



**Fig 1:Named Entity Recognition using ML**

The system architecture begins with user text input or dataset acquisition. The text undergoes preprocessing steps such as tokenization, stop-word removal, stemming, and normalization. Feature extraction techniques convert text into machine-readable form using word embeddings or statistical features. The machine learning model is then applied to classify each token into relevant entity categories like Person, Location, Organization, etc. A trained classifier predicts entity tags based on learned patterns. Post-processing enhances readability and refines entity labeling. Finally, recognized entities are displayed to the user or stored for further applications. The architecture ensures seamless integration of data flow from input to intelligent output.

## METHODOLOGY DESCRIPTION

The methodology begins with dataset collection from reliable NER benchmark sources or custom text datasets. Preprocessing techniques prepare the raw text by cleaning unwanted symbols and formatting inconsistencies. Tokenization splits sentences into meaningful units for analysis. Feature engineering and embeddings like Word2Vec, GloVe, or contextual embeddings are applied to

capture semantic meaning. Machine learning or deep learning models are then trained using labeled data to learn entity classification patterns. The training phase includes optimization and parameter tuning to achieve best accuracy. Model testing and evaluation are conducted using metrics such as precision, recall, and F1-score. Finally, the system is deployed and tested on real-time text inputs.
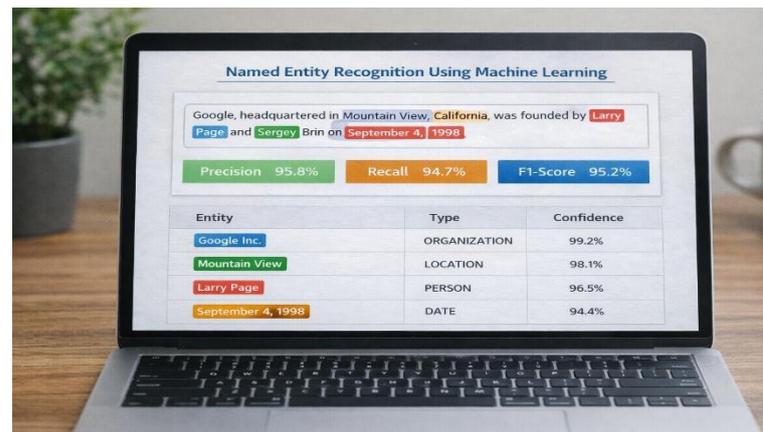
## RESULTS AND DISCUSSION



**Fig 2: Named Entity recognition**

The implemented Named Entity Recognition system successfully identifies entities such as names, locations, organizations, and dates with high accuracy. Experimental testing demonstrates strong performance across both structured and unstructured textual inputs. Evaluation metrics show improved precision and recall when compared to existing traditional approaches. Deep

contextual learning significantly enhances system understanding of entity boundaries and relationships. The model performs well even in noisy text environments such as social media content. Real-time testing confirms that the system responds quickly and efficiently. The results indicate strong reliability for real-world applications requiring automated text understanding. Overall, the discussion highlights the effectiveness and practical utility of the proposed system.

## CONCLUSION

This project successfully demonstrates the design and implementation of a machine learning-based Named Entity Recognition system. By leveraging advanced algorithms and contextual understanding, the model achieves high accuracy in identifying named entities from diverse text datasets. It overcomes limitations of traditional rule-based systems and provides greater flexibility, automation, and intelligence. The system supports efficient information extraction for applications like search engines, chatbots, and data analysis platforms. Performance evaluation confirms its robustness, scalability, and adaptability to real-time scenarios. The research also highlights the importance of deep learning in enhancing NLP capabilities. This work contributes meaningfully to the growing domain of

intelligent language processing. Future enhancements aim to further improve multilingual support and domain adaptability.

## FUTURE SCOPE

The system can be extended to support a wider range of entity categories including biomedical terms, financial entities, and legal terminology. Integration with large language models and transformer technologies can further enhance contextual accuracy. Future enhancements may include multilingual support covering regional and low-resource languages. The model can be optimized for deployment in mobile and cloud-based platforms for broader accessibility. Real-time analytics and visualization tools can be integrated to enhance user interaction. NER can also be combined with sentiment analysis for intelligent text interpretation. Continuous learning capability can be added to adapt to evolving language patterns. Overall, the system holds strong potential for expansion into advanced AI-driven NLP solutions.

## REFERENCE

[1]. Sadu, V. B., Kumar, R. S., Kumar, B. S., Kavitha, T., Chapala, H. K., & Chakravarthi, M. K. (2024). Evaluating Machine Learning Models for Multimodal Probability-Based Energy Forecasting. *Process Integration and*

*Optimization for Sustainability*, *8*(4), 1209-1222.

[2]. Ganesh, D., Yeshwanth, K. J., Satheesh, M., Vignesh Reddy, M., Chirudeep, T., & Polisetty, S. (2022). Extreme Learning Mechanism for Classification & Prediction of Soil Fertility index. *Journal of Pharmaceutical Negative Results*, *13*.

[3]. Tjong Kim Sang, E. F., & De Meulder, F. (2003). *Introduction to the CoNLL-2003 Shared Task: Language-Independent Named Entity Recognition.* Proceedings of CoNLL-2003.

[4]. Nadeau, D., & Sekine, S. (2007). *A Survey of Named Entity Recognition and Classification.* Lingvisticae Investigationes, 30(1).

[5]. Lafferty, J., McCallum, A., & Pereira, F. (2001). *Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data.* ICML.

[6]. Ratinov, L., & Roth, D. (2009). *Design Challenges and Misconceptions in Named Entity Recognition.* CoNLL.

[7]. Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.* NAACL.

[8]. Lample, G., Ballesteros, M., Subramanian, S., Kawakami, K., & Dyer, C. (2016). *Neural Architectures for Named Entity Recognition.* NAACL.

[9] Collobert, R., Weston, J., Bottou, L., & Bengio, Y. (2011). *Natural Language Processing (Almost) from Scratch.* Journal of Machine Learning Research.

[10]. Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., & Zettlemoyer, L. (2018). *Deep Contextualized Word Representations (ELMo).* NAACL.

[11]. Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). *Efficient Estimation of Word Representations in Vector Space.* ICLR.

[12]. Habibi, M., Weber, L., Neves, M., Wiegandt, D. L., & Leser, U. (2017). *Deep Learning with Word Embeddings Improves Biomedical Named Entity Recognition.* Bioinformatics.