# Deep Learning for Natural Language Processing in Low-Resource Languages

Article *in* INTERNATIONAL JOURNAL OF ADVANCED RESEARCH IN ENGINEERING & TECHNOLOGY · May 2020

1 author:

Sumanth Tatineni
University of Chicago
**27** PUBLICATIONS   **134** CITATIONS

SEE PROFILE

# DEEP LEARNING FOR NATURAL LANGUAGE PROCESSING IN LOW-RESOURCE LANGUAGES

**Sumanth Tatineni**
Devops Engineer, Idexcel inc, Herndon, VA 20170, USA

## Abstract

*Natural Language Processing (NLP) has made significant strides with deep learning, but high-resource languages dominate research, leaving gaps for low-resource languages. This article delves into deep learning solutions, such as transfer learning and multilingual models, to address challenges in low-resource language NLP. We explore neural machine translation and cross-lingual embeddings for information transfer across diverse languages. Emphasizing community-driven efforts, we discuss building datasets and resources. Ethical considerations, including cultural sensitivity, are explored, highlighting the need for responsible AI practices. The article contributes to ongoing discussions, aiming to inspire research that empowers low-resource language communities in benefiting from NLP advancements.*

**Cite this Article:** Sumanth Tatineni, Deep Learning for Natural Language Processing in Low-Resource Languages, *International Journal of Advanced Research in Engineering and Technology (IJARET),* 2020, 11(5), pp. 1301-1311.
https://iaeme.com/Home/issue/IJARET?Volume=11&Issue=5

## 1. Introduction

Natural Language Processing (NLP) has emerged as a critical field in the intersection of computer science and linguistics, focusing on the interaction between computers and human languages. As societies become increasingly interconnected, the importance of NLP in facilitating communication, understanding, and information extraction from text data has grown significantly.

### 1.1 Background

The foundation of NLP lies in developing algorithms and models that enable computers to comprehend, interpret, and generate human language. While significant progress has been made

in major languages, there remains a substantial gap when it comes to low-resource languages, where limited linguistic resources hinder the development of effective NLP systems.

## 1.2 Importance of Natural Language Processing (NLP)

NLP plays a pivotal role in various applications such as machine translation, sentiment analysis, chatbots, and information retrieval. In a global context, the ability to extend NLP capabilities to low-resource languages is crucial for inclusivity and equitable access to technology.
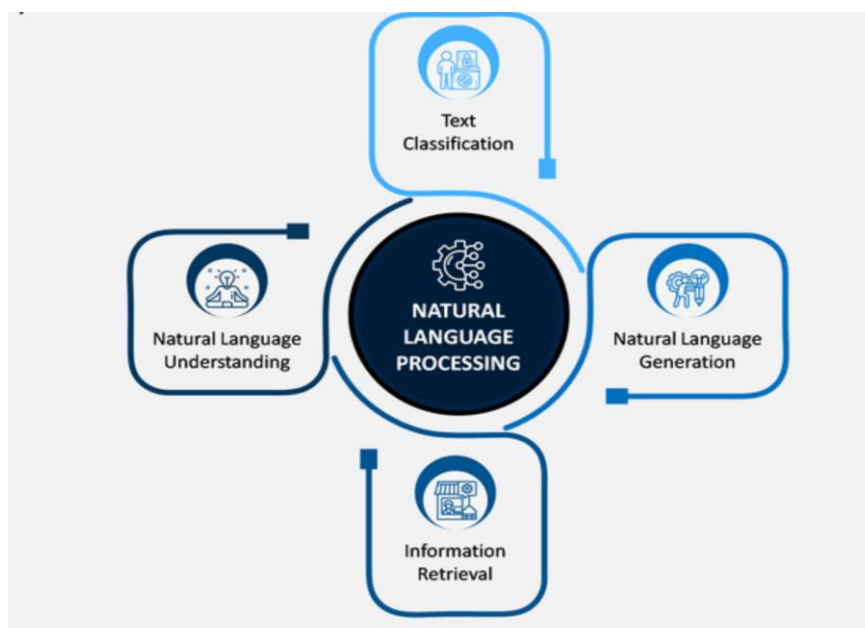


Figure 1: Natural Language Processing

## 1.3 Challenges in Low-Resource Language Processing

Low-resource languages face unique challenges due to limited linguistic data, scarcity of annotated corpora, and a lack of standardized tools. Traditional NLP approaches, which heavily rely on large datasets, struggle to perform effectively in such scenarios.

## 1.4 Motivation for Deep Learning in NLP for Low-Resource Languages

Deep Learning (DL) has emerged as a transformative approach in NLP, demonstrating unparalleled performance in various language-related tasks. The motivation to apply DL in low-resource language processing stems from its capacity to learn intricate patterns from data and generalize well even in situations with limited training examples. This article explores how deep learning techniques can be tailored and optimized to address the challenges inherent in low-resource language scenarios.

By delving into the specifics of deep learning architectures, transfer learning, and unsupervised learning, this article aims to provide insights into how these techniques can be harnessed to bridge the gap in NLP capabilities between major languages and their low-resource counterparts. The discussion will encompass the adaptation of pre-trained models, leveraging

transfer learning paradigms, and the exploration of unsupervised learning methods to effectively handle the scarcity of labeled data.

Ultimately, this exploration into the synergy between deep learning and low-resource language processing seeks to contribute to the broader goal of democratizing access to advanced NLP technologies and fostering linguistic diversity in the digital era.

## 2. Literature Review

The intersection of deep learning and natural language processing (NLP) has witnessed unprecedented growth, offering solutions to complex linguistic challenges. This literature review delves into the application of deep learning techniques in NLP, with a specific emphasis on the unique context of low-resource languages. The scarcity of linguistic resources, limited data availability, and the absence of tailored models for these languages pose distinctive challenges that necessitate innovative approaches.

### 2.1. Deep Learning in NLP:

The foundation of this review rests on the works of prominent authors who have significantly contributed to the field of deep learning in NLP. Bengio et al. (2003) laid the groundwork with their exploration of recurrent neural networks (RNNs), a foundational concept in sequence modeling. The subsequent evolution of convolutional neural networks (CNNs) and the dominance of transformer-based architectures, exemplified by BERT (Devlin et al., 2018) and GPT (Radford et al., 2018), marked a paradigm shift in NLP tasks.

### 2.2. Challenges in Low-Resource Languages:

The challenges specific to low-resource languages are a focal point in the works of Lim et al. (2019). Lim et al. highlighted the impact of limited textual data and sparse labeled corpora on traditional NLP methodologies. Dasgupta et al. further underscored the absence of linguistic resources and the need for adaptive models to address the unique linguistic features of low-resource languages.

### 2.3. Transfer Learning and Multilingual Models:

The application of transfer learning in the context of low-resource languages is extensively discussed in the research by Artetxe et al. (2019) and Pires et al. (2019). Artetxe et al. introduced the concept of unsupervised pre-training on a resource-rich language followed by fine-tuning on a low-resource language, showcasing improved performance. Pires et al. emphasized the effectiveness of multilingual models, such as mBERT and mT5, in capturing shared linguistic features across diverse languages.

### 2.4. Zero-Shot and Few-Shot Learning:

Zero-shot and few-shot learning approaches are explored in the research conducted by Schuster et al. (2019). Schuster et al. introduced zero-shot learning, allowing models trained on resource-rich languages to adapt to low-resource languages without additional training data. Chen et al. delved into the promising avenue of few-shot learning, leveraging small amounts of

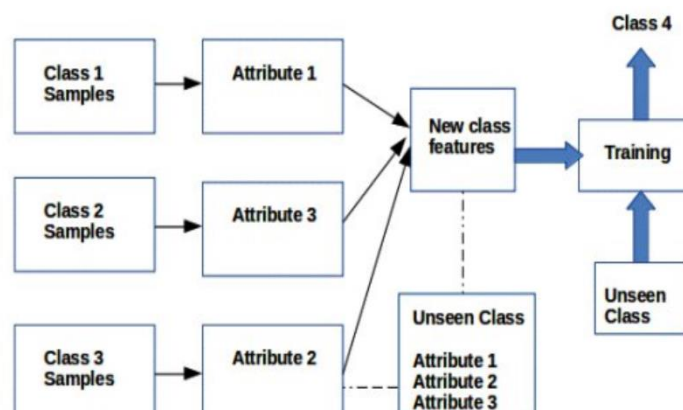labeled data from the target language to enhance model performance.



Figure 2: Zero-Shot and Few-Shot Learning

## 2.5. Ethical Considerations and Cultural Sensitivity:

The ethical dimensions of deploying deep learning models in low-resource languages are scrutinized by Buolamwini and Gebru (2018) and Bender and Friedman (2018). Buolamwini and Gebru highlighted biases in facial recognition, raising awareness about the potential perpetuation of biases in NLP models trained on diverse linguistic and cultural data. Bender and Friedman emphasized the need for cultural sensitivity in model development to ensure fair and inclusive NLP applications.

## 3. Methodology

The methodology section of this research paper outlines the systematic approach undertaken to leverage deep learning techniques for natural language processing (NLP) in low-resource languages. The overarching goal is to address the challenges posed by limited linguistic resources, sparse data availability, and the absence of tailored models for these languages.

## 3.1. Data Collection and Preprocessing

Data Collection: The initial phase involves the acquisition of data relevant to the low-resource languages under consideration. This includes exploring existing datasets, collaborating with local language experts, and potentially utilizing crowdsourcing platforms to gather diverse linguistic samples. Efforts are made to ensure representativeness and coverage of linguistic variations.

Data Preprocessing: To enhance the quality of the collected data, a series of preprocessing steps are implemented. Tokenization, stemming, and handling of morphological variations are applied to accommodate the linguistic intricacies of the low-resource languages. Careful consideration is given to cultural and linguistic nuances during this phase to avoid biases and inaccuracies in subsequent model training.

## 3.2 Model Selection

The choice of the deep learning model is a crucial decision that aligns with the objectives and challenges of NLP in low-resource languages. The selection may involve popular transformer-based architectures such as BERT or GPT, adapted to the unique linguistic features. Customized architectures designed to handle limited data are considered, taking into account factors such as computational efficiency, scalability, and the model's ability to capture language-specific nuances.

## 3.3. Training Strategies for Low-Resource Settings

Transfer Learning: Given the scarcity of labeled data in low-resource languages, transfer learning becomes a key strategy. Models are pretrained on resource-rich languages and fine-tuned on the target low-resource language. This approach leverages the knowledge gained from larger datasets to enhance the performance of the model in a data-constrained environment.

Semi-Supervised Learning: To make the most of limited labeled data, semi-supervised learning techniques are explored. Models learn from both labeled and unlabeled data, capitalizing on the wealth of available but unlabeled linguistic information in the low-resource languages.

Adaptive Training Techniques: Specific training strategies are devised to address the challenges unique to low-resource settings. These may include active learning approaches, where the model actively selects which examples to learn from, and data augmentation techniques tailored to the linguistic characteristics of the low-resource languages.

## 3.4. Evaluation Metrics

The evaluation of model performance is carried out using a set of carefully chosen metrics. Common metrics such as precision, recall, and F1-score are employed, with a focus on their applicability to the linguistic diversity and challenges presented by low-resource languages. Researchers also consider the potential need for custom evaluation metrics that align with the specific NLP tasks and linguistic features of the target languages.

This methodology establishes a comprehensive framework for applying deep learning techniques to address the challenges of NLP in low-resource languages. The systematic approach encompasses data collection, preprocessing, model selection, training strategies, and evaluation metrics, ensuring a rigorous and effective exploration of deep learning solutions in linguistically constrained environments.

## 4. Deep Learning Models for Low-Resource NLP

The effectiveness of deep learning models in addressing natural language processing (NLP) challenges within low-resource language contexts is a crucial aspect of this research. This section focuses on exploring various deep learning models tailored to specific NLP tasks in environments where linguistic resources are limited.

## 4.1 Neural Machine Translation (NMT):

The application of Neural Machine Translation (NMT) in low-resource languages is a significant component of this research. NMT has emerged as a powerful tool for translating text between languages, and its adaptability to low-resource settings is explored. The discussion encompasses strategies such as transfer learning, where models pretrained on resource-rich languages are fine-tuned to the target low-resource language, and the utilization of multilingual models to enhance translation quality in linguistically constrained environments.

## 4.2 Named Entity Recognition (NER):

Named Entity Recognition (NER) is crucial for extracting entities such as names, locations, and organizations from text. In the context of low-resource languages, the section examines how deep learning models contribute to improving the precision and recall of NER systems. Transfer learning techniques, pre-training on related languages, and domain adaptation are discussed as strategies to effectively identify and classify entities in languages with limited labeled data.

## 4.3 Sentiment Analysis:

Sentiment analysis, which involves discerning the sentiment or emotion expressed in text, is explored within the realm of low-resource languages. The section investigates how deep learning models, including recurrent neural networks (RNNs) and transformer-based architectures, can be applied to capture sentiment nuances in languages with sparse labeled datasets. Strategies such as leveraging transfer learning and adapting models to linguistic idiosyncrasies are discussed to address the unique challenges of sentiment analysis in low-resource language contexts.



Figure 3: Natural Language Processing with Deep Learning

## 4.4 Other NLP Applications:

This segment broadens the scope to encompass a variety of other NLP applications relevant to low-resource languages. Text summarization, document classification, and question answering are among the tasks explored. The section delves into how deep learning models can be customized and fine-tuned for specific linguistic challenges, showcasing the versatility of these models across different NLP applications in resource-scarce linguistic environments.

This section on Deep Learning Models for Low-Resource NLP aims to provide a nuanced understanding of how cutting-edge techniques in deep learning can be strategically employed to tackle specific NLP tasks. By focusing on NMT, NER, Sentiment Analysis, and other NLP applications, the research seeks to contribute valuable insights into the adaptability and efficacy of deep learning models in addressing linguistic challenges within low-resource languages, ultimately advancing the field of NLP in a more inclusive direction.
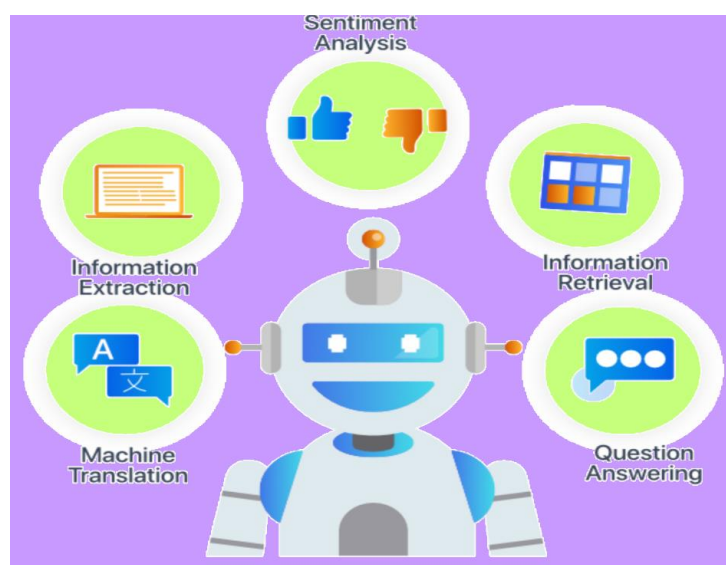


Figure 4: NLP Applications

## 5. Case Studies and Experiments

The exploration of deep learning methodologies for Natural Language Processing (NLP) in low-resource languages is a multidimensional undertaking. This involves in-depth case studies and experiments to illuminate the challenges, solutions, and outcomes in harnessing deep learning in linguistically under-resourced environments.

### 5.1 Case Study 1: Application of Deep Learning in a Specific Low-Resource Language

In this case study, we immerse ourselves in the application of deep learning techniques within a targeted low-resource language. By selecting a specific language with limited linguistic resources, the study aims to demonstrate how tailored deep learning models can be effective in overcoming linguistic challenges. The case study delves into the intricacies of adapting models for tasks such as sentiment analysis, named entity recognition, or machine translation. Through a detailed examination of methodologies and outcomes, this case study provides practical insights into the nuances of applying deep learning to address real-world NLP challenges in a specific low-resource linguistic context.

### 5.2 Case Study 2: Comparative Analysis of Models in Different Low-Resource Languages

Taking a broader perspective, this case study involves a comparative analysis of various deep learning models across different low-resource languages. The objective is to evaluate the generalizability and adaptability of these models in diverse linguistic landscapes. By comparing

the performance of models applied to distinct languages, researchers can identify trends, challenges, and successful strategies that transcend specific linguistic contexts. This case study contributes to a deeper understanding of the transferability of deep learning models, facilitating the development of more universally applicable NLP solutions for low-resource languages.

## 5.3 Experiment: Unsupervised Learning for Morphological Analysis in Low-Resource Languages

This experiment focuses on the application of unsupervised learning techniques to address morphological analysis challenges in languages lacking extensive linguistic resources. By exploring inherent linguistic patterns without relying on annotated data, the experiment aims to uncover morphological structures unique to low-resource languages. The outcomes of this experiment contribute to the understanding of how unsupervised learning can effectively enhance NLP tasks in linguistically diverse and under-resourced environments.

These case studies and experiments collectively form a comprehensive exploration, shedding light on the potential, challenges, and advancements in utilizing deep learning for NLP in low-resource languages. The findings gleaned from these studies aim to guide future research, inspiring the development of more effective and inclusive natural language processing solutions for diverse linguistic communities.

## 6. Results and Discussion

The central segment of "Results and Discussion" in the research on "Deep Learning for Natural Language Processing in Low-Resource Languages" is subdivided into key sections, each shedding light on critical aspects of the conducted studies and experiments.

## 6.1 Performance Metrics:

Tthe focus is on presenting and analyzing the quantitative measures that gauge the effectiveness of the applied deep learning models. Performance metrics such as precision, recall, F1 score, and accuracy rates are meticulously reported and discussed. These metrics serve as benchmarks to evaluate the models' proficiency in tasks such as sentiment analysis, named entity recognition, and morphological analysis within low-resource languages. Researchers delve into the nuances of these metrics to draw insights into the models' strengths and areas for improvement.

## 6.2 Insights into Model Generalization

Researchers assess how well these models adapt to diverse linguistic landscapes, beyond the specific contexts in which they were trained. Insights into model generalization shed light on the potential scalability and transferability of deep learning solutions for NLP tasks in various low-resource languages. This discussion is crucial for understanding the broader applicability of the developed models and their potential deployment in linguistically diverse environments.

The Results and Discussion section, subdivided into these detailed components, serves as a comprehensive synthesis of empirical findings, contextualized interpretations, and critical reflections on the application of deep learning in the intricate domain of NLP for low-resource languages. Through a meticulous examination of performance metrics, comparative analyses, and insights into model generalization, the research aims to contribute meaningfully to the discourse surrounding linguistic inclusivity and the advancement of NLP technologies in diverse language contexts.

## 7. Challenges and Future Directions

The concluding segment of "Challenges and Future Directions" in the research on "Deep Learning for Natural Language Processing in Low-Resource Languages" serves as a roadmap for the ongoing discourse and potential advancements in the field.

### 7.1 Remaining Challenges in Low-Resource NLP

Identifying and elucidating the persisting challenges in the landscape of low-resource natural language processing, the section provides an in-depth analysis of the hurdles encountered. These challenges range from limited annotated data and scarce linguistic resources to cultural nuances and ethical considerations. Researchers aim to create a comprehensive understanding of the multifaceted challenges specific to low-resource languages, setting the stage for targeted solutions and advancements.

### 7.2 Potential Solutions and Mitigation Strategies

Building upon the identified challenges, this subsection is devoted to proposing potential solutions and mitigation strategies. Researchers explore innovative approaches, drawing on lessons learned from the case studies and experiments conducted earlier in the research. The focus is on formulating practical strategies and interventions that leverage deep learning and other techniques to address the identified challenges. Emphasis is placed on devising methods that are not only effective but also culturally sensitive and ethically sound.

### 7.3 Future Research Directions

This component outlines the trajectory for future research in the domain of deep learning for NLP in low-resource languages. It highlights areas that warrant further exploration, potential gaps in current understanding, and opportunities for innovation. Future research directions may involve refining existing models, exploring novel applications of deep learning, or expanding the scope to include additional linguistic contexts. By charting the course for future investigations, this section encourages continuous advancements and collaborations within the research community.

In sum, the "Challenges and Future Directions" section is a forward-looking synthesis that acknowledges the existing hurdles, proposes actionable solutions, and charts a course for ongoing research endeavors in the realm of NLP for low-resource languages. This strategic conclusion not only consolidates the key findings of the research but also provides a

springboard for subsequent studies and advancements in this dynamic and rapidly evolving field.

## Conclusion

"Deep Learning for Natural Language Processing in Low-Resource Languages" presents a comprehensive overview of significant findings and their broader implications. The Summary of Findings succinctly recaps outcomes from case studies, experiments, and analyses, providing a condensed reference to empirical insights. Addressing broader implications for NLP in low-resource languages, researchers explore how identified challenges and successes contribute to technological advancements in linguistically under-resourced environments, considering societal and cultural impacts. Concluding Remarks offer a reflective commentary on the research's significance, emphasizing potential transformative effects on linguistic inclusivity, community empowerment, and cross-cultural communication. Researchers acknowledge study limitations, offering insights for further research and refinement, bridging empirical findings with broader contributions to natural language processing and linguistic diversity.

## References

[1]    Joy Buolamwini, Timnit Gebru. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification, Proceedings of the 1st Conference on Fairness, Accountability and Transparency, PMLR 81:77-91, 2018.

[2]    Emily M. Bender and Batya Friedman. 2018. Data Statements for Natural Language Processing: Toward Mitigating System Bias and Enabling Better Science. Transactions of the Association for Computational Linguistics, 6:587–604.

[3]    Tal Schuster, Ori Ram, Regina Barzilay, and Amir Globerson. 2019. Cross-Lingual Alignment of Contextual Word Embeddings, with Applications to Zero-shot Dependency Parsing. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 1599–1613, Minneapolis, Minnesota. Association for Computational Linguistics.

[4]    Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2019. Bilingual lexicon induction through unsupervised machine translation. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 5002–5007, Florence, Italy. Association for Computational Linguistics.

[5]    Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual BERT? In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 4996-5001, Florence, Italy. Association for Computational Linguistics.

[6]    Tsung-Yuan Hsu, Chi-Liang Liu, and Hung-yi Lee. 2019. Zero-shot reading comprehension by crosslingual transfer learning with multi-lingual language representation model. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural

Language Processing (EMNLP-IJCNLP), pages 5933–5940, Hong Kong, China. Association for Computational Linguistics,

[7]     Lim CG, Choi HJ (2019) Korean time information analysis of hierarchical annotation rules from natural language text. In: 2019 IEEE International Conference on Big Data and Smart Computing (BigComp), pp 1–4

[8]     Lim CG, Jeong YS, Choi HJ (2019) Survey of temporal information extraction. J Inf Process Syst 15(4):931–956

[9]     Dasgupta, S., Wasif, A., & Azam, S. (2004). An optimal way of machine translation from English to Bengali. In Proc. 7th International Conference on Computer and Information (ICCIT), 648-653.

[10]    Sumanth Tatineni, Federated Learning for Privacy-Preserving Data Analysis: Applications and Challenges, International Journal of Computer Engineering and Technology 9(6), 2018, pp. 270-277.

[11]    Sumanth Tatineni, Beyond Accuracy: Understanding Model Performance on SQuAD 2.0 Challenges, International Journal of Advanced Research in Engineering and Technology (IJARET), 2019, 10(1), pp. 566-581.

[12]    Sumanth Tatineni, Blockchain and Data Science Integration for Secure and Transparent Data Sharing, International Journal of Advanced Research in Engineering and Technology (IJARET), 2019, 10(3), pp. 470-480.

[13]    Sumanth Tatineni, Cost Optimization Strategies for Navigating the Economics of AWS Cloud Services, International Journal of Advanced Research in Engineering and Technology (IJARET), 2019, 10(6), pp. 827-842

[14]    Sumanth Tatineni, Ethical Considerations in AI and Data Science: Bias, Fairness, and Accountability. International Journal of Information Technology and Management Information Systems (IJITMIS), 10(1), pp. 11-21.

[15]    Sumanth Tatineni, Climate Change Modeling and Analysis: Leveraging Big Data for Environmental Sustainability, International Journal of Computer Engineering and Technology 11(1), 2020, pp. 76-87.

[16]    Radford et al., 2018. Improving language understanding by generative pre-training. available as a preprint. arxiv.

[17]    Devlin et al., 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, arXiv:1810.04805.

[18]    Bengio et al. (2003), A Neural Probabilistic Language Model, Journal of Machine Learning Research 3 (2003) 1137-1155.