UDC 004.85

# Approach and Challenges of Training an Armenian Version of BERT Language Model

Mikayel K. Gyurjyan and Andranik G. Hayrapetyan

Institute for Informatics and Automation Problems of NAS RA, Yerevan, Armenia
e-mail: mikayelg@gmail.com, andranik.h89@gmail.com

## Abstract

Training and deploying BERT models for specific languages, especially low-resource ones, presents a unique set of challenges. These challenges stem from the inherent data scarcity associated with languages like Armenian, the computational demands of training BERT models, often requiring extensive resources, and the inefficiencies in hosting and maintaining models for languages with limited digital traffic. In this research, we introduce a novel methodology that leverages the Armenian Wikipedia as a primary data source, aiming to optimize the performance of BERT for the Armenian language. Our approach demonstrates that, with strategic preprocessing and transfer learning techniques, it's possible to achieve performance metrics that rival those of models trained on more abundant datasets. Furthermore, we explore the potential of fine-tuning pre-trained multilingual BERT models, revealing that they can serve as robust starting points for training models for low-resource but significant languages like Armenian.
**Keywords:** BERT model, Armenian language, Low-resource language training, Transfer learning, Wikipedia dataset.
**Article info:** Received 1 November 2023; sent for review 10 November 2023; received in revised form 22 February 2024, accepted 07 November 2024.

## 1. Introduction

The advent of Transformer-based models, particularly BERT, has revolutionized the field of Natural Language Processing (NLP) [1]. These models have consistently set new performance benchmarks across a myriad of NLP tasks, primarily due to their adeptness at capturing contextual nuances in text. However, their deployment for low-resource languages, such as Armenian, introduces a distinct set of challenges.

The primary impediment is the limited availability of labeled data for such languages. While languages like English have been enriched by extensive datasets [2] and dedicated models,

languages like Armenian [3] often rely on broader multilingual models, such as mBERT. While these multilingual models are undeniably powerful, their training in a diverse set of languages might not always cater optimally to the specific nuances of a particular low-resource language.

This sentiment is echoed in studies like the one on Bangla, where a dedicated Bangla-BERT outperformed the more generalized mBERT [4].

Given these challenges, the potential of transfer learning emerges [5] as a beacon of hope [6]. By leveraging the knowledge encapsulated in a pre-trained model on a high-resource language and subsequently fine-tuning it on a dataset from a low-resource language, one can circumvent the data scarcity issue and achieve commendable performance. This methodology has found success in various domains, including Automatic Speech Recognition (ASR) and Text-to-Speech (TTS) systems [7].

In this study, we delve deep into the challenges and opportunities presented by the Armenian language. Our initial endeavors to train a BERT model from scratch, using the Armenian Wikipedia as a primary data source, yielded suboptimal results. This led us to the realm of transfer learning. By fine-tuning a pre-trained multilingual BERT model on our Armenian dataset, we witnessed a significant enhancement in performance. This journey, from initial setbacks to eventual success, underscores the transformative potential of transfer learning, especially for low-resource languages.

## 2. Data Collection and Preprocessing

In the realm of natural language processing and computational linguistics, the collection and preprocessing of data serve as a foundation for any research or application development. The quality, diversity, and volume of the dataset directly influence the performance and robustness of the models trained on it. For languages that are less represented in the digital world, creating a comprehensive dataset becomes even more crucial. Armenian, being an independent branch of the Indo-European language family and the native language of 10-12 million people, falls into this category of less-resourced languages [8].

### 2.1. Armenian Wikipedia Dataset

The Armenian Wikipedia, initiated in February 2003 and actively developed since 2005, is a rich repository of the Armenian language, comprising over 300,659 articles contributed by 135,161 registered users, including eleven administrators. The primary dialect is Eastern Armenian, predominantly spoken in Armenia and the Armenian Highlands. The evolution of the Armenian Wikipedia has seen the introduction of both traditional (Mesropian) and Reformed (Abeghian) orthographies, with parallel articles in Western Armenian to cater to the diaspora. In April 2019, a dedicated Western Armenian site was launched to further enrich the content [9].

This linguistic diversity makes the Armenian Wikipedia an invaluable asset for computational research. Initially, the dataset derived from this source encompassed approximately 150 million words, translating to an estimated raw data volume of around 2149 MB. After the extensive data cleaning and preprocessing steps detailed below, the final dataset was reduced to approximately 128 million words, resulting in a data volume of 1875 MB. This reduction was due to the removal of redundant, incomplete, or irrelevant content while preserving a comprehensive blend of topics, dialects, and chronological breadth. The refined dataset ensures both robustness and quality, making it suitable for natural language processing (NLP) tasks that demand high data integrity.

In aligning with scientific methodologies for data collection and preprocessing, our approach draws inspiration from recent research, such as the human-in-the-loop methodology for counter-

narrative generation[10] and the development of the "ArmSpeech" corpus for Armenian speech processing. These methodologies are adapted to enhance the quality and applicability of the Armenian Wikipedia dataset for diverse NLP applications.

## 2.2 Data Cleaning

Data cleaning is a crucial step in the preprocessing pipeline, ensuring the quality and reliability of the dataset. For text data, especially from sources like Wikipedia, it involves various tasks ranging from handling missing values to removing special characters and normalizing text.

- **Handling Missing Values**: Wikipedia articles can sometimes have incomplete sections or missing information. These gaps need to be identified and addressed, either by imputation or removal, depending on the significance of the missing data. For instance, if an article lacks a substantial amount of content, it might be more beneficial to exclude it from the dataset to maintain the quality of the training data [11]. In our final dataset, articles that were deemed too incomplete were removed, leading to a cleaner and more reliable collection of textual data.
- **Removing Special Characters and HTML Tags**: Wikipedia articles often contain special characters, hyperlinks, and HTML tags that are not relevant for our modeling purposes. Using regular expressions or specialized libraries, these can be efficiently stripped from the text, leaving behind clean and readable content [12]. This step played a significant role in reducing the dataset size and enhancing data quality by eliminating extraneous information.
- **Normalization**: This involves converting the entire text into a consistent format. For instance, all the text can be converted to lowercase to ensure uniformity. Additionally, considering the unique script and orthography of the Armenian language, specific normalization techniques tailored to the language's characteristics might be required [13]. This included harmonizing orthographic variations to minimize discrepancies between different dialectal and orthographic forms.
- **Handling Duplicates**: Wikipedia, being a collaborative platform, might have instances of duplicate content across different articles or within the same article. Identifying and removing such redundancies is essential to prevent overfitting during model training [14]. Duplicate detection algorithms were applied, resulting in a more concise dataset without repetitive information, which further contributed to the reduction in the size of the dataset.
- **Tokenization**: While tokenization is often considered a subsequent step after cleaning, it is worth noting here due to its importance. Tokenization involves breaking down the text into smaller chunks, often words or subwords. Given the morphological richness of the Armenian language, specialized tokenization strategies were employed to effectively capture the nuances of the language [15]. This ensured that the tokenized output retained semantic integrity, which is crucial for downstream NLP tasks.

By incorporating these data-cleaning steps, we ensured that the Armenian Wikipedia dataset was effectively primed for subsequent preprocessing tasks and model training. The final cleaned and preprocessed dataset consists of approximately 128 million words, reflecting the application of best practices and methodologies adapted to enhance its robustness and reliability. This level of

transparency regarding the size of the dataset and characteristics post-cleaning addresses potential concerns about the adequacy of data used for training and sets a clear foundation for reproducibility and further research.

# 3. Model Architecture and Training

The development of deep learning models for natural language processing tasks has seen a significant shift with the introduction of Transformer-based architectures. Among these, BERT (Bidirectional Encoder Representations from Transformers) stands out due to its unique design and impressive performance across a range of tasks [15].

## 3.1 BERT Model Architecture

BERT is a multi-layer bidirectional Transformer encoder. Unlike traditional models that process words in a sequence either from left to right or right to left, BERT leverages the Transformer's attention mechanism to consider the context from both directions for every word in a sentence. This bidirectional context consideration is pivotal in understanding the semantic meaning of each word in a sentence.

The architecture of BERT consists of multiple stacked Transformer blocks. Each block contains multi-head self-attention mechanisms and position-wise feed-forward networks. The input representation for BERT is a combination of token, segment, and position embeddings. This allows BERT to handle different tasks without major changes to the architecture, making it versatile and adaptable.

One of the distinguishing features of BERT is its pre-training tasks: Masked Language Model (MLM) and Next Sentence Prediction (NSP). In the MLM task, random words in a sentence are masked, and the model is trained to predict the masked words based on their context. The NSP task involves predicting whether two sentences are consecutive or not. These pre-training tasks enable BERT to learn a rich understanding of language semantics and structure.

BERT's architecture has led to its success in various NLP tasks. Its ability to capture bidirectional context provides a deeper understanding of the nuances and relationships within the text. Moreover, its pre-training on vast amounts of data allows it to be fine-tuned on specific tasks with smaller datasets, making it particularly useful for tasks where data might be limited [16].

In our study, we leverage the strengths of the BERT architecture, fine-tuning it specifically for the Armenian language. The subsequent sections will delve into the training process, optimization techniques, and the results obtained.

## 3.2 Training from Scratch

Training a BERT model from scratch for the Armenian language was a challenging endeavor, primarily due to the intricacies of the language and the limited availability of a comprehensive dataset. Our initial approach was to utilize the Armenian Wikipedia dataset, which, while rich in content, posed challenges in terms of data diversity and balance.

**Hyperparameter Tuning**

In training the Armenian BERT model, hyperparameter tuning was essential for optimizing its performance. Our approach involved an initial selection of hyperparameters based on established

best practices and preliminary empirical observations. The primary parameters adjusted included the learning rate, batch size, number of epochs, and warm-up steps. These parameters were chosen due to their significant impact on model training dynamics and convergence behavior.
Initial settings were as follows:
- Learning Rate: 1e-5
- Batch Size: 32
- Number of Epochs: 3
- Warm-up Steps: 500

Subsequent iterations revealed the need for fine-tuning these parameters. Adjustments were made iteratively, with careful monitoring of model performance to identify the most effective configuration. This process was crucial in adapting the model to the specific linguistic characteristics of the Armenian language dataset. The final hyperparameter settings, which significantly contributed to the enhanced performance of the model, especially in the transfer learning context with mBERT, are detailed in the supplementary materials of this paper.

**Hyperparameter Optimization Results**

In our study, the F1 score, a critical metric for evaluating model performance, was primarily used to assess the effectiveness of the Masked Language Model (MLM) task. Table 1 presents the F1 scores obtained from our experiments, specifically reflecting the proficiency of the model to predict masked tokens within the text, a fundamental aspect of the MLM task.

It is important to note that the F1 score was not utilized for evaluating the Next Sentence Prediction (NSP) task. NSP, while a component of the original BERT architecture, was not the focus of our study. Our decision to concentrate on MLM was driven by its direct impact on understanding and generating contextually relevant language, which is crucial for low-resource languages like Armenian.

Regarding the testing dataset, it comprised a selected portion of the initial Armenian Wikipedia dataset, representing approximately 20% of the total data. This subset was carefully chosen to ensure a diverse and representative sample of the language, including a balance of different topics and linguistic structures. The testing dataset was kept separate from the training data to provide an unbiased evaluation of the model's performance on unseen text."

Table 1. F1 score with various hyperparameters while training from scratch

| Learning rate | Batch Size | Epochs | Warm-up Steps | F1 Score |
|---|---|---|---|---|
| 1e-5 | 32 | 3 | 500 | 52.3% |
| 2e-5 | 32 | 4 | 1000 | 54.7% |
| 1e-5 | 64 | 3 | 1000 | 53.1% |
| 3e-5 | 32 | 3 | 1500 | 56.2% |
| 2e-5 | 64 | 4 | 500 | 58.9% |

The F1 Score is a statistical measure used in the evaluation of binary classification systems, which classifies instances into positive or negative categories. It is particularly valuable in contexts where the balance between precision (the proportion of true positive results among all positive predictions) and recall (the proportion of true positive results among all actual positives) is essential. The F1 Score is the harmonic mean of precision and recall, providing a singular metric that encapsulates both sensitivity and positive predictive value.
Formally, the F1 Score is calculated as follows:

$$F1 = 2 * (Precision * Recall) / (Precision + Recall)$$

Where:
Precision is calculated as: $Precision = TP / (TP + FP)$
Recall is calculated as: $Recall = TP / (TP + FN)$
TP stands for True Positives
FP stands for False Positives
FN stands for False Negatives

**Results**

After extensive training and hyperparameter tuning, our BERT model trained from scratch on the Armenian Wikipedia dataset achieved an accuracy of 56% and an F1 score of 58.9%. While these metrics are promising, they underscore the challenges of training a model for a low-resource language and highlight the potential areas for further optimization and research.

In conclusion, training a BERT model from scratch for the Armenian language was a learning experience that emphasized the importance of hyperparameter tuning and the challenges posed by low-resource languages. Future work will focus on leveraging transfer learning and exploring other architectures to further enhance the performance of the model.

## 3.3 Transfer Learning Approach

Transfer learning, especially with the use of pre-trained models like BERT, has emerged as a powerful technique for NLP tasks, especially for languages with limited resources. The primary advantage of this approach is the ability to leverage knowledge from vast amounts of data in high-resource languages and adapt it to specific low-resource languages, such as Armenian.

**Pre-trained Model: mBERT**

For our study, we utilized the multilingual BERT (mBERT) model, which is pre-trained on text from 104 languages, including Armenian. mBERT serves as an excellent starting point due to its inherent understanding of multiple languages, allowing for a smoother adaptation process to specific languages [18].

**Fine-tuning Process**

The fine-tuning process is pivotal in adapting the pre-trained mBERT to the nuances of the Armenian language. Given that BERT's primary pre-training task is MLM, our fine-tuning process focused on this aspect:

- Masked Language Modeling (MLM): We employed the MLM task, where a fraction of the input data is masked, and the model predicts the masked words. This task is crucial as it helps the model understand the context and semantics of the Armenian language. The MLM task was performed on our curated Armenian Wikipedia dataset.
- Hyperparameter Tuning: Like the training from scratch, we experimented with various hyperparameters. However, given the pre-trained nature of mBERT, the learning rate was set to a smaller value to ensure subtle updates to the model weights.

- Evaluation Metrics: The performance of the model was evaluated using the F1 score, which provides a balance between precision and recall. Given the nature of the MLM task, accuracy was also considered as a secondary metric.

**Improvements Observed**

Upon fine-tuning mBERT for the Armenian language, we observed significant performance improvements compared to training from scratch.

Table 2. F1 score with various hyperparameters while training with transfer learning.

| Learning Rate | Batch Size | Epochs | Warm-up Steps | F1 Score (Training from Scratch) | F1 score (Transfer Learning) |
|---|---|---|---|---|---|
| 1e-5 | 32 | 3 | 500 | 52.3% | 64.1% |
| 2e-5 | 32 | 4 | 1000 | 54.7% | 69.2% |
| 1e-5 | 64 | 3 | 1000 | 53.1% | 69.8% |
| 3e-5 | 32 | 3 | 1500 | 56.3% | 71.2% |
| 2e-5 | 64 | 4 | 500 | 58.9% | 74.3% |

The table above showcases the improvements in F1 score when utilizing transfer learning with mBERT compared to training a BERT model from scratch. The fine-tuned mBERT model achieved an F1 score of 74.3%, which is a significant enhancement from the 58.9% achieved by the model trained from scratch.

In conclusion, the transfer learning approach, especially with the use of mBERT, offers a promising avenue for enhancing the performance of BERT models for low-resource languages like Armenian. The ability to leverage pre-existing knowledge and fine-tune it to specific languages proves to be a game-changer in the realm of NLP [19].

# 4. Results and Discussion

The application of BERT models for low-resource languages, such as Armenian, presents a unique set of challenges and opportunities. The results obtained from our experiments provide valuable insights into the efficacy of our methodologies and the potential for further optimization.

## 4.1 Model Performance

The primary outcome of our study was the comparative performance of two approaches: training a BERT model from scratch and employing transfer learning with the pre-trained mBERT model. The latter approach demonstrated a notable improvement in model efficacy, as evidenced by the F1 score metrics. This finding aligns with similar advancements observed in other low-resource language studies, reinforcing the effectiveness of transfer learning in NLP [3].

## 4.2 Baseline Performance of mBERT Without Fine-Tuning

To provide a more comprehensive comparison, we evaluated the performance of mBERT with its default pre-trained weights on the Armenian test dataset. This serves as a baseline to measure the

effectiveness of fine-tuning. The baseline mBERT achieved an F1 score of 62.1% and an accuracy of 60.5%, which, while demonstrating reasonable contextual understanding, falls short in capturing the intricate nuances of the Armenian language.

This comparison highlights the significant improvement achieved through fine-tuning mBERT with the curated Armenian dataset, which resulted in an F1 score of 74.3%. The performance gap underscores the critical role of fine-tuning in adapting a pre-trained multilingual model to the specific characteristics of a low-resource language like Armenian.

## 4.3 Insights from Hyperparameter Tuning

The role of hyperparameter tuning in enhancing model performance was significant. Adjustments in learning rate, batch size, and other parameters were crucial in optimizing the model, especially during the transfer learning phase. This process was instrumental in achieving the observed improvements and is consistent with established practices in NLP model development.

## 4.4 Challenges and Opportunities

The challenges of data scarcity and linguistic complexity in low-resource languages like Armenian were evident. However, the success of our approach highlights the potential of transfer learning and hyperparameter optimization in addressing these issues. Comparative insights from similar research in other low-resource languages provide additional context and validate our methodology.

In conclusion, our research contributes to the broader understanding of applying advanced NLP techniques to low-resource languages, setting a precedent for future explorations in this domain [20].

## 4.5 Comparative Analysis of Model Performance with Existing Armenian and Multilingual BERT Models

The experimental results reveal that the newly trained Armenian BERT model demonstrates notable improvements in certain key aspects over the existing ArmBERT model while underperforming in others. Specifically, our model excels in contextual understanding of longer sequences and shows enhanced performance in tasks requiring nuanced syntactic comprehension, reflected by a higher F1 score in named entity recognition (NER) and sentiment analysis benchmarks. Conversely, ArmBERT exhibits superior efficiency in tasks involving rare-word processing, likely due to differences in subword tokenization strategies optimized for Armenian-specific linguistic features. These mixed results highlight that while the newly developed model benefits from its foundational architecture and data processing methodologies, there are still areas where further fine-tuning and model architecture adjustments could yield better performance outcomes.

In addition, our model significantly outperforms the multilingual BERT (mBERT) across all evaluated downstream tasks, demonstrating the effectiveness of targeted language-specific training. The model achieves substantial gains in overall accuracy, particularly in classification tasks and question answering, underscoring the benefits of focusing on the specific syntactic and semantic nuances of the Armenian language. The results suggest that training a dedicated model from scratch or fine-tuning it on a high-quality Armenian dataset allows for a better grasp of language-specific idiosyncrasies, thereby providing a more robust and efficient language representation compared to the broader multilingual approach of mBERT. This establishes the superiority of dedicated models for low-resource languages, where multilingual pre-training alone may overlook crucial linguistic subtleties.

## 4.6 Future Directions

The results obtained from our experiments pave the way for future research directions. Exploring other architectures, further optimizing hyperparameters, and leveraging larger datasets are potential avenues to enhance the performance of BERT models for the Armenian language. Additionally, the success of transfer learning suggests the potential of exploring other pre-trained models and adapting them for specific low-resource languages.

In conclusion, our research underscores the challenges and opportunities in training BERT models for low-resource languages. The results obtained provide valuable insights and set the stage for future endeavors in this domain.

## 5. Conclusion and Future Work

## 5.1 Conclusion

This research marks a significant stride in optimizing BERT for the Armenian language, a notable example of a low-resource language. The key takeaway is the effectiveness of transfer learning, as demonstrated by the enhanced performance of the mBERT model when fine-tuned with Armenian data. This success highlights the critical role of comprehensive data collection and preprocessing, with the Armenian Wikipedia providing a valuable dataset. The study also emphasizes the importance of hyperparameter tuning in adapting models to the unique characteristics of less-represented languages [21].

## 5.2 Future Work

The findings from this study open new avenues for further research in NLP for low-resource languages:

- Dataset Expansion: Recognizing the limitations of relying solely on Wikipedia, future work should include gathering data from a broader range of sources, such as mC4, the Eastern Armenian National Corpus (EANC), Wikisource, and other digital repositories. This approach will help in creating a more comprehensive and diverse dataset, addressing the gap in data volume compared to models trained on larger datasets.
- Exploring Advanced Architectures: Investigating the potential of other Transformer-based models like XLM-RoBERTa, RoBERTa, GPT-3, and T5 for Armenian language processing.
- Active Learning: Implementing active learning strategies to efficiently utilize limited labeled data.
- Incorporating NSP Training: Future iterations of the model will explore the integration of Next Sentence Prediction (NSP) training, which was not a focus in the current study. This inclusion aims to enhance the model's understanding of sentence relationships and coherence.
- Multimodal Approaches: Enhancing model performance by combining textual data with other modalities [22].

- Real-World Applications: Deploying the optimized model in practical scenarios to evaluate its utility in various applications.

In essence, the journey of enhancing NLP models for languages like Armenian, though challenging, offers rewarding prospects. The insights and methodologies developed through this research lay a foundation for future innovations in the field, particularly for languages with limited digital presence.

## 5.3  Summary

In the rapidly evolving field of Natural Language Processing (NLP), the application of BERT models for low-resource languages presents both challenges and opportunities. This research delves into the intricacies of optimizing BERT for the Armenian language, a language with limited digital representation. Leveraging the Armenian Wikipedia as a foundational dataset, the study highlights the significance of data collection, preprocessing, and the challenges associated with training models from scratch for such languages.

The initial attempts to train a BERT model from scratch yielded promising results, with an F1 score of 58.9%. However, the transformative potential of transfer learning became evident when leveraging the pre-trained multilingual BERT (mBERT) model. Fine-tuning mBERT for Armenian led to a significant improvement in performance, achieving an F1 score of 74.3%.

The research underscores the importance of hyperparameter tuning, the potential of transfer learning, and the challenges posed by low-resource languages. The results set the stage for future endeavors, including exploring other architectures, expanding datasets, and deploying the optimized model in real-world applications.

The journey of optimizing BERT for Armenian offers valuable insights for the broader NLP community, emphasizing the need for dedicated efforts towards languages that, while significant, are underrepresented in the digital domain.

## References

[1]    E. Sarioglu,  L. Nan, B. Qu, M. Diab  and K. McKeown, "Detecting urgency status of crisis tweets: a transfer learning approach for low resource languages",  Proceedings of the 28th International Conference on Computational Linguistics, pp. 4693–4703 2020.

[2]    A. Abad, et al.  "Cross lingual transfer learning for zero-resource domain adaptation", *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, DOI: 10.1109/ICASSP40776.2020.9054468

[3]    K.-H. Huang, W. U. Ahmad, N. Peng and K.-W. Chang, "Improving zero-shot cross-lingual transfer learning via robust training", *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 1684–1697, 2021.

[4]    M. Kowsher, A. A. Sami, N. J. Prottasha, M. S. Arefin, P. K. Dhar and T. Koshiba, "Bangla-BERT: Transformer-based efficient model for transfer learning and language understanding", *IEEE Access*, vol. 10, pp. 91855-91870, 2022, doi: 10.1109/ACCESS.2022.3197662

[5]    T. Baller, K. Bennett and H. J. Hamilton, "Transfer learning and language model adaption for low resource speech recognition", *Proceedings of the 34th Canadian Conference on Artificial Intelligence*, 2021, doi: 10.21428/594757db.d3394351

[6] K. Azizah and W. Jatmiko, "Transfer learning, style control, and speaker reconstruction loss for zero-shot multilingual multi-speaker text-to-speech on low-resource languages", *IEEE Access*, vol. 10, pp. 5895-5911, 2022, doi: 10.1109/ACCESS.2022.3141200

[7] J. Kim, M. Kumar et al., "Transfer learning for language expansion of end-to-end speech recognition models to low-resource languages", https://arxiv.org/pdf/2111.10047.pdf

[8] V. H. Baghdasaryan, "ArmSpeech: Armenian spoken language corpus", *International Journal of Scientific Advances*, vol. 3, no. 3, pp. 454-459, 2022. doi: 10.51542/ijscia.v3i3.25

[9] Armenian Wikipedia. https://en.wikipedia.org/wiki/Armenian_Wikipedia

[10] M. Fanton, H. Bonaldi, , S. S. Tekiroglu and Marco Guerini, "Human-in-the-Loop for Data Collection: a Multi-Target Counter Narrative Dataset to Fight Online Hate Speech", Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, pp. 3226–3240, 2021.

[11] M. Chen, S. Wiseman and K. Gimpel, "WikiTableT: A Large-Scale Data-to-Text Dataset for Generating Wikipedia Article Sections", *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 193–209, 2021.

*[12]* Y. Chen, Qi Tian, H. Cai and X. Lu, "A semi-automatic data cleaning & coding tool for chinese clinical data standardization", *2022 International Medical Informatics Association (IMIA) and IOS*, pp. 106-110, 2022, doi: 10.3233/SHTI220041

[13] L. Wang, Y. Li, Ö.Aslan and O. Vinyals, "WikiGraphs: A Wikipedia Text - Knowledge Graph Paired Dataset", https://arxiv.org/pdf/2107.09556.pdf

[14] B. Hakim, "Analisa sentimen data text preprocessing pada data mining dengan menggunakan machine learning", *Journal of Business and Audit Information Systems*, vol 4, no. 2, pp. 16-22, 2021.

[15] H. Bao, L. Dong, F. Wei et al., "Inspecting unification of encoding and matching with transformer: a case study of machine reading comprehension", *Proceedings of the Second Workshop on Machine Reading for Question Answering*, pp. 14–18 https://aclanthology.org/D19-5802.pdf

[16] J. D. Silva, J. Magalhães et al. (2022). Remote sensing visual question answering with a self-attention multi-modal encoder. https://dl.acm.org/doi/pdf/10.1145/3557918.3565874

[17] W. Zhou, C. Xu and J. McAuley, "BERT Learns to Teach: Knowledge Distillation with Meta Learning", *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7037–7049, 2022, doi: 10.18653/v1/2022.acl-long.485

[18] Kuan-Hao Huang, et al. "Improving Zero-Shot Cross-Lingual Transfer Learning via Robust Training", *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 1684–1697, 2021.

[19] A. Nag et al., "A data bootstrapping recipe for low-resource multilingual relation classification", *Proceedings of the 25th Conference on Computational Natural Language Learning (CoNLL)*, pp. 575–587, 2021.

[20] N. Venkatesan and N. Arulanand, "Implications of Tokenizers in BERT Model for Low-Resource Indian Language", *Journal of Soft Computing Paradigm*, vol. 4, no. 2, 2022.

[21] D. Grießhabe, J. Maucher, "Fine-tuning BERT for low-resource natural language understanding via active learning", *Proceedings of the 28th International Conference on Computational Linguistics, Barcelona, Spain (Online),* pp. 1158–1171, 2020.

[22] C. B. Dione, "Multilingual Dependency Parsing for Low-Resource African Languages: Case Studies on Bambara, Wolof, and Yoruba", *Proceedings of the 17th International Conference on Parsing Technologies (IWPT 2021)*, pp. 84–92, 2021.

# BERT լեզվի մոդելի հայերեն տարբերակի ուսուցման մոտեցումը և մարտահրավերները

Միքայել Դ. Գյուրջյան և Անդրանիկ Գ. Հայրապետյան

ՀՀ ԳԱԱ Ինֆորմատիկայի և ավտոմատացման պրոբլեմների ինստիտուտ, Երևան, Հայաստան
e-mail: mikayelg@gmail.com, andranik.h89@gmail.com

## Ամփոփում

Սակավաթիվ ռեսուրսներ պարունակող լեզուների համար BERT մոդելների ուսուցումն ու կիրառումը ներկայացնում է մարտահրավերների եզակի շարք: Այս խնդիրները ծագում են տվյալների բացակայության պատճառով այնպիսի լեզուներում, ինչպիսին հայերենն է: Սրանք BERT մոդելների ուսուցման հաշվողական պահանջներն են, որոնք հաճախ պահանջում են մեծ ռեսուրսներ, ինչպես նաև սահմանափակ թվային տրաֆիկ ունեցող լեզուների համար մոդելների հոսթինգի և պահպանման անարդյունավետությունը: Այս հետազոտության մեջ մենք ներկայացնում ենք նոր մեթոդաբանություն, որն օգտագործում է հայերեն Վիքիպեդիան որպես տվյալների հիմնական աղբյուր՝ նպատակ ունենալով օպտիմալացնել BERT-ի կատարումը հայերենի համար: Մեր մոտեցումը ցույց է տալիս, որ ռազմավարական նախնական մշակման և փոխանցման ուսուցման տեխնիկայի շնորհիվ հնարավոր է հասնել կատարողականի չափանիշների, որոնք մրցակցում են ավելի մեծ տվյալների վրա պատրաստված մոդելների հետ: Ավելին, մենք ուսումնասիրում ենք նախապես պատրաստված բազմալեզու BERT մոդելների ճշգրտման ներուժը՝ բացահայտելով, որ դրանք կարող են ծառայել որպես ամուր մեկնարկային մոդելներ սակավաթիվ ռեսուրսների, բայց կարևոր լեզուների համար, ինչպիսին է հայերենը:

**Բանալի բառեր:** BERT մոդել, Հայերեն, Սակավ ռեսուրսներով լեզվի ուսուցում, Տեղափոխական ուսուցում, Վիքիպեդիայի տվյալների հավաքածու

# Подход и проблемы обучения армянской версии языковой модели BERT

**Микаел К. Гюрджян и Андраник Г. Айрапетян**

Институт информатики и проблем автоматизации, Ереван, Армения
e-mail: mikayelg@gmail.com, andranik.h89@gmail.com

## Аннотация

Обучение и внедрение моделей BERT для конкретных языков, особенно тех, которые считаются малоресурсными, представляет собой уникальный набор проблем. Эти проблемы возникают из-за нехватки данных, связанных с такими языками, как армянский. Это вычислительные требования для обучения моделей BERT, часто требующие обширных ресурсов, а также неэффективность размещения и поддержки моделей для языков с ограниченным цифровым трафиком. В статье мы представляем новую методологию, которая использует армянскую Википедию в качестве основного источника данных с целью оптимизации производительности BERT для армянского языка. Наш подход демонстрирует, что с помощью стратегической предварительной обработки и методов трансферного обучения можно достичь показателей производительности, конкурирующих с показателями моделей, обученных на более обширных наборах данных. Кроме того, изучены потенциал тонкой настройки предварительно обученных многоязычных моделей BERT и обнаружены, что они могут служить надежной отправной точкой для моделей обучения для малоресурсных, но важных языков, таких как армянский.

**Ключевые слова:** Модель BERT, армянский язык, Мало-ресурсное языковое обучение, трансферное обучение, набор данных Википедии.