
PRE-TRAINING POLISH TRANSFORMER-BASED LANGUAGE MODELS AT SCALE

A PREPRINT

Sławomir Dadas

National Information Processing Institute
Warsaw, Poland
sdadas@opi.org.pl

Michał Perełkiewicz

National Information Processing Institute
Warsaw, Poland
mperelkiewicz@opi.org.pl

Rafał Poświata

National Information Processing Institute
Warsaw, Poland
rposwiata@opi.org.pl

June 11, 2020

ABSTRACT

Transformer-based language models are now widely used in Natural Language Processing (NLP). This statement is especially true for English language, in which many pre-trained models utilizing transformer-based architecture have been published in recent years. This has driven forward the state of the art for a variety of standard NLP tasks such as classification, regression, and sequence labeling, as well as text-to-text tasks, such as machine translation, question answering, or summarization. The situation have been different for low-resource languages, such as Polish, however. Although some transformer-based language models for Polish are available, none of them have come close to the scale, in terms of corpus size and the number of parameters, of the largest English-language models. In this study, we present two language models for Polish based on the popular BERT architecture. The larger model was trained on a dataset consisting of over 1 billion polish sentences, or 135GB of raw text. We describe our methodology for collecting the data, preparing the corpus, and pre-training the model. We then evaluate our models on thirteen Polish linguistic tasks, and demonstrate improvements over previous approaches in eleven of them.

Keywords Language Modeling · Natural Language Processing

1 Introduction

Unsupervised pre-training for Natural Language Processing (NLP) has gained popularity in recent years. The goal of this approach is to train a model on a large corpus of unlabeled text, and then use the representations the model generates as an input for downstream linguistic tasks. The initial popularization of these methods was related to the successful applications of pre-trained word vectors (embeddings), the most notable of which include Word2Vec [31], GloVe [35], and FastText [5]. These representations have contributed greatly to the development of NLP. However, one of the main drawbacks of such tools was that the static word vectors did not encode contextual information. The problem was addressed in later studies by proposing context-dependent representations of words based on pre-trained neural language models. For this purpose, several language model architectures which utilize bidirectional long short-term memory (LSTM) layers have been introduced. The popular models such as ELMo [36], ULMFiT [18], and Flair [1], have led to significant improvements in a wide variety of linguistic tasks. Shortly after, Devlin et al. [16] introduced BERT - a different type of language model based on transformer [43] architecture. Instead of predicting the next word in a sequence, BERT is trained to reconstruct the original sentence from one in which some tokens have

been replaced by a special *mask token*. Since the text representations generated by BERT have proved to be effective for NLP problems - even those which were previously considered challenging, such as question answering or common sense reasoning - more focus has been put on transformer-based language models. As a result, in the last two years we have seen a number of new methods based on that idea, with some modifications in the architecture or the training objectives. The approaches that have gained wide recognition include RoBERTa [27], Transformer-XL [13], XLNet [47], Albert [25], and Reformer [20].

The vast majority of research on both transformer-based language models and transfer learning for NLP is targeted toward the English language. This progress does not translate easily to other languages. In order to benefit from recent advancements, language-specific research communities must adapt and replicate studies conducted in English to their native languages. Unfortunately, the cost of training state-of-the-art language models is growing rapidly [34], which makes not only individual scientists, but also some research institutions unable to reproduce experiments in their own languages. Therefore, we believe that it is particularly important to share the results of research - especially pre-trained models, datasets, and source code of the experiments - for the benefit of the whole scientific community. In this article, we describe our methodology for training two language models for Polish language based on BERT architecture. The smaller model follows the hyperparameters of an English-language BERT-base model, and the larger version follows the BERT-large model. To the best of our knowledge, the latter is the largest language model for Polish available to date, both in terms of the number of parameters (355M) and the size of the training corpus (135GB). We have released both pre-trained models publicly¹. We evaluate our models on several linguistic tasks in Polish, including nine from the KLEJ benchmark [39], and four additional tasks. The evaluation covers a set of typical NLP problems, such as binary and multi-class classification, textual entailment, semantic relatedness, ranking, and Named Entity Recognition (NER).

1.1 Language-specific and multilingual transformer-based models

In this section we provide an overview of models based on the transformer architecture for languages other than English. Apart from English, the language on which NLP research is most focused currently is Chinese. This is reflected in the number of pre-trained models available [46, 9, 16, 42]. Other languages for which we found publicly available pre-trained models included: Arabic [3], Dutch [14, 15], Finnish [44], French [30, 26], German, Greek, Italian, Japanese, Korean, Malaysian, Polish, Portuguese [41], Russian [24], Spanish [6], Swedish, Turkish, and Vietnamese [32]. Models covering a few languages of the same family are also available, such as SlavicBERT (Bulgarian, Czech, Polish, and Russian) [4] and NordicBERT² (Danish, Norwegian, Swedish, and Finnish). The topic of massive multilingual models covering tens, or in some cases more than a hundred languages, has attracted more attention in recent years. The original BERT model [16] was released along with a multilingual version covering 104 languages. XLM [7] (fifteen, seventeen and 100 languages) and XLM-R [8] (100 languages) were released in 2019. Although it was possible to use these models for languages in which no monolingual models were available, language-specific pre-training usually leads to better performance. To date, two BERT-base models have been made available for Polish: HerBERT [39] and Polbert³, both of which utilize BERT-base architecture.

1.2 Contributions

Our contributions are as follows: 1) We trained two transformer-based language models for Polish, consistent with the BERT-base and BERT-large architectures. To the best of our knowledge, the second model is the largest language model trained for Polish to date, both in terms of the number of parameters and the size of the training corpus. 2) We proposed a method for collecting and pre-processing the data from the Common Crawl database to obtain clean, high-quality text corpora. 3) We conducted a comprehensive evaluation of our models on thirteen Polish linguistic tasks, comparing them to other available transformer-based models, as well as recent state-of-the-art approaches. 4) We made the source code of our experiments available to the public, along with the pre-trained models.

2 Language model pre-training

In this section, we describe our methodology for collecting and pre-processing the data used for training BERT-base language models. We then present the details of the training, explaining our procedure and the selection of hyperparameters used in both models.

¹<https://github.com/sdadas/polish-roberta>

²https://github.com/botxo/nordic_bert

³<https://github.com/kldarek/polbert>

2.1 Training corpus

Transformer-based models are known for their high capacity [19, 22], which means that they can benefit from large quantities of text. An important step in the process of creating a language model, therefore, is to collect a sufficiently large text corpus. We have taken into account that the quality of the text used for training will also affect the final performance of the model. The easiest way to collect a large language-specific corpus is to extract it from Common Crawl - a public web archive containing petabytes of data crawled from web pages. The difficulty with this approach is that web-based data is often noisy and unrepresentative of typical language use, which could eventually have a negative impact on the quality of the model. In response to this, we have developed a procedure for filtering and cleaning the Common Crawl data to obtain a high-quality web corpus. The procedure is as follows:

1. We download full HTML pages (WARC files in Common Crawl), and use the resulting metadata to filter the documents written in Polish language.
2. We use *Newspaper3k*⁴ - a tool which implements a number of heuristics for extracting the *main content* of the page, discarding any other text such as headers, footers, advertisements, menus, or user comments.
3. We then remove all texts shorter than 100 characters. Additionally, we identify documents containing the words: ‘przeglądarka’, ‘ciasteczka’, ‘cookies’, or ‘javascript’. The presence of these words may indicate that the extracted content is a description of a cookie policy, or default content for browsers without JavaScript enabled. We discard all such texts if they are shorter than 500 characters.
4. In the next step, we use a simple statistical language model (KenLM [17]), trained on a small Polish language corpus to assess the quality of each extracted document. For each text, we compute the perplexity value and discard all texts with perplexity higher than 1000.
5. Finally, we remove all duplicated texts.

The full training corpus we collected is approximately 135GB in size, and is composed of two components: the web part and the base part. For the web part, which amounts to 115GB of the corpus, we downloaded three monthly dumps of Common Crawl data, from November 2019 to January 2020, and followed the pre-processing steps described above. The base part, which comprises the remaining 20GB, is composed of publicly available Polish text corpora: the Polish language version of Wikipedia (1.5GB), the Polish Parliamentary Corpus (5GB), and a number of smaller corpora from the CLARIN (<http://clarin-pl.eu>) and OPUS (<http://opus.nlpl.eu>) projects, as well as Polish books and articles.

2.2 Training procedure

The authors of the original BERT paper [16] proposed two versions of their transformer-based language model: BERT-large (more parameters and higher computational cost), and BERT-base (fewer parameters, more computationally efficient). To train the models for Polish language, we adapted the same architectures. Let L denote the number of encoder blocks, H denote the hidden size of the token representation, and A denote the number of attention heads. Specifically, we used $L = 12, H = 768, A = 12$ for the base model, and $L = 24, H = 1024, A = 16$ for the large model. The large model was trained on the full 135GB text corpus, and the base model on only the 20GB base part. The training procedure we employed is similar to the one suggested in the RoBERTa pre-training approach [27]. Originally, BERT utilized two training objectives - Masked Language Modeling (MLM), and Next Sentence Prediction (NSP). We trained our models with the MLM objective, since it has been shown that NSP fails to improve the performance of the pre-trained models on downstream tasks [27]. We also used dynamic token masking, and trained the model with a larger batch size than the original BERT. The base model was trained with a batch size of 8000 sequences for 125 000 training steps: the large model was trained with a batch size of 30 000 sequences for 50 000 steps. The reason for using such a large batch size for the bigger model is to stabilize the training process. During our experiments, we observed significant variations in training loss for smaller batch sizes, indicating that the initial combination of learning rate and batch size had caused an exploding gradient problem. To address the issue, we increased the batch size until the loss stabilized.

Both models were pre-trained with the Adam optimizer using the following optimization hyperparameters: $\epsilon = 1e-6, \beta_1 = 0.9, \beta_2 = 0.98$. We utilized a learning rate scheduler with linear decay. The learning rate is first increased for a warm-up phase of 10 000 update steps to reach a peak of $7e-4$, and then linearly decreased for the remainder of the training. We also mimicked the dropout approach of the original BERT model: a dropout of 0.1 is applied on all layers and attention weights. The maximum length of a sequence was set to 512 tokens. We do not combine sentences from the training corpus: each is treated as a separate training sample. To encode input sequences

⁴<https://newspaper.readthedocs.io/en/latest/>

into tokens, we employed SentencePiece [23] Byte Pair Encoding (BPE) algorithm, and set the maximum vocabulary size to 50 000 tokens.

3 Evaluation

In this section, we discuss the process and results of evaluating our language models on thirteen Polish downstream tasks. Nine of these tasks constitute the recently developed KLEJ benchmark [39]; three of them have already been introduced in Dadas et al. [12]; and the last, named entity recognition, was a part of the PolEval⁵ evaluation challenge. First, we compare the performance of our models with other Polish and multilingual language models evaluated on the KLEJ benchmark. Next, we present detailed per-task results, comparing our models with the previous state-of-the-art solutions for each of the tasks.

3.1 Task descriptions

NKJP (The National Corpus of Polish (Narodowy Korpus Języka Polskiego)) [37] is one of the largest text corpora of the Polish language, consisting of texts from Polish books, news articles, web content, and transcriptions of spoken conversations. A part of the corpus, known as the ‘one million subcorpus’, contains annotations of named entities from six categories: ‘persName’, ‘orgName’, ‘geogName’, ‘placeName’, ‘date’, and ‘time’. The authors of the KLEJ benchmark used this subset to create a named entity classification task [39]. First, they filtered out all sentences containing entities of more than one type. Next, they randomly assigned sentences to train development and test sets according to the rule that each named entity mentioned appears in only one of these splits. They undersample the ‘persName’ class, and merge the ‘date’ and ‘time’ classes to increase class balance. Finally, they selected sentences without any named entity, and assigned them the ‘noEntity’ label. The resulting dataset consisted of 20 000 sentences belonging to six classes. The task is to predict the presence and type of each named entity. Classification accuracy is also reported.

STAGS is a corpus created by Dadas et al. [12] for their study on the subject of sentence representations in Polish language. This dataset was created automatically by extracting sentences from headlines and short descriptions of articles posted on the Polish social network, wykop.pl. It contains approximately 50 000 sentences, all longer than thirty characters, from eight popular categories: film, history, food, medicine, automotive, work, sport, and technology. The task is to assign a sentence to one of these classes in which classification accuracy is the measure.

CBD (Cyberbullying Detection) [38] is a binary classification task, the goal of which is to determine whether a Twitter message constitutes a case of cyberbullying or not. This was a sub-task of task 6 in the PolEval 2019 competition. The dataset prepared by the competition’s organizers contains 11 041 tweets, extracted from nineteen of the most popular Polish Twitter accounts in 2017. The F1-score was used to measure the performance of the models.

DYK ‘Did you know?’ (*‘Czy wiesz?’*) [28] is a dataset used for the evaluation and development of Polish language question answering systems. It consists of 4721 question-answer pairs obtained from the *Czy wiesz...* Polish Wikipedia project. The answer to each question was found in the linked Wikipedia article. Rybak et al. [39] used this dataset to devise a binary classification task, the goal of which is to predict whether the answer to the given question is correct or not [39]. Positive responses were additionally marked within larger fragments of responded text. Negative samples were selected by the BPE token overlap between a question and a possible answer. The F1-score was also reported for this task.

PSC The Polish Summaries Corpus [33] is a corpus of manually created summaries of Polish language news articles. The dataset contains both abstract free-word summaries and extraction-based summaries created by selecting text spans from the original documents. Based on PSC, [39] formulated a text-similarity task [39]. They generate positive pairs by matching each extractive summary with the two least similar abstractive ones in the same article. Negative pairs were obtained by finding the two most similar abstractive summaries for each extractive summary, but from different articles. To calculate the similarity between summaries, they used the BPE token overlap. The F1-score was used for evaluation.

PolEmo2.0 [21] is a corpus of consumer reviews obtained from four domains: medicine, hotels, products, and school. Each of the reviews is annotated with one of four labels: positive, negative, neutral, or ambiguous. In general, the task is to choose the correct label, although here two special versions of the task are distinguished: PolEmo2.0-IN and PolEmo2.0-OUT. In PolEmo2.0-IN, both the training and test sets come from the same domains, namely medicine and hotels. In PolEmo2.0-OUT, however, the test set comes from the product and school domains. In both cases, accuracy was used for evaluation.

⁵<http://2018.poleval.pl/index.php/tasks>

Allegro Reviews (AR) [39] is a sentiment analysis dataset of product reviews from the e-commerce marketplace, allegro.pl. Each review has a rating on a five-point scale, in which one is negative, and five is positive. The task is to predict the rating of a given review. The macro-average of the mean absolute error per class (wMAE) is applied for evaluation.

CDSC (The Compositional Distributional Semantics Corpus) [45] is a corpus of 10 000 human-annotated sentence pairs for semantic relatedness and entailment, in which image captions from forty-six thematic groups were used as sentences. Two tasks are proposed based on this dataset. The CDSC-R problem involves predicting the relatedness between a pair of sentences, on a scale of zero to five, in which zero indicates that the sentences are not related, and five indicates that they are highly related. In this task, the Spearman correlation is used as an evaluation measure. CDSC-E’s task is to classify whether the premise entails the hypothesis (entailment), negates the hypothesis (contradiction), or is unrelated (neutral). For this task, accuracy is reported.

SICK [12] is a manually translated Polish language version of the English Natural Language Inference (NLI) corpus, SICK (Sentences Involving Compositional Knowledge) [29], and consists of 10 000 sentence pairs. As with the CDSC dataset, two tasks can also be distinguished here. SICK-R is the task of predicting the probability distribution of relatedness scores (ranging from 1 to 5) for the sentence pair, in which the Spearman correlation is used for evaluation. SICK-E is a multiclass classification problem in which the relationship between two sentences is classified as entailment, contradiction, or neutral. Accuracy is used once again to measure performance.

PolEval-NER 2018 [2] was task 2 in the PolEval 2018 competition, the goal of which was to detect and assign the correct category and subcategory (if applicable) to a found named entity. In this study the task was simplified, as only the main categories had to be found. The effectiveness of the models is verified by the F1-score measure. This task was prepared on the basis of the NKJP dataset previously presented.

3.2 Task-specific fine-tuning

To evaluate our language models on downstream tasks, we fine-tuned them separately for each task. In our experiments, we encounter three types of problem: classification, regression, and Named Entity Recognition (NER). In classification tasks, the model is expected to predict a label from a set of two or more classes. Regression concerns the prediction of a continuous numerical value. NER is a special case of sequence tagging, i.e. predicting a label for each element in a sequence. The dataset for each problem consists of training and test parts, and in most cases also includes a validation part. The general fine-tuning procedure is as follows: we train our model on the training part of the dataset for a specific number of epochs. If the validation set is available, we compute the validation loss after each epoch, and select the model checkpoint with the best validation loss. For datasets without a validation set, we select the last epoch checkpoint. Then, we perform an evaluation on the test set using the selected checkpoint.

In the case of classification and regression tasks, we attach an additional fully-connected layer to the output of the $[CLS]$ token, which always remains in the first position of a sequence. For classification, the number of outputs for this layer is equal to the number of classes, and the softmax activation function is used. For regression, it is a linear layer with a single output. The models are fine-tuned with the Adam optimizer using the following hyperparameters: $\epsilon = 1e-6$, $\beta_1 = 0.9$, $\beta_2 = 0.98$. A learning rate scheduler with polynomial decay is utilized. The first 6% of the training steps are reserved for the warm-up phase, in which the learning rate is gradually increased to reach a peak of $1e-5$. By default, we train for ten epochs with a batch size of sixteen sequences. The specific fine-tuning steps and exceptions to the procedure are discussed below:

- **Classification on imbalanced datasets** – Some of the binary classification datasets considered in the evaluation, such as CBD, DYK, and PSC, are imbalanced, which means that they contain significantly fewer samples of the first class than of the second class. To counter this imbalance, we utilize a simple resampling technique: samples for the minority class in the training set are duplicated, and some samples for the majority class are randomly discarded. We set the resampling factor to 3 for the minority class, and 1 (DYK, PSC) or 0.75 (CBD) respectively for the majority class. Additionally, we increase the batch size for those tasks to thirty-two.
- **Regression** - In many cases, a regression task is restricted to a specific range of values for which the prediction is valid. For example, Allegro Reviews contains user reviews with ratings between one and five stars. For fine-tuning, we scale all the outputs of regression models to be within the range of $[0, 1]$, and then rescale them to their original range during evaluation. Before rescaling, any negative prediction is set to 0, and any prediction greater than 1 is limited to 1.
- **Named entity recognition** - Since sequence tagging, in which the model is expected to generate per-token predictions, is different from simple classification or regression tasks, we decided to adapt an existing named entity recognition approach for fine-tuning using our language models. For this purpose, we employed a method from Shibuya and Hovy [40], who proposed a transformer-based named entity recognition model with a Conditional Ran-

Table 1: Results on the KLEJ benchmark.

Model	Average	NKJP	CDSC-E	CDSC-R	CBD	PE2-I	PE2-O	DYK	PSC	AR
Base models										
mBERT	79.5	91.4	93.8	92.9	40.0	85.0	66.6	64.2	97.9	83.3
SlavicBERT	79.8	93.3	93.7	93.3	43.1	87.1	67.6	57.4	98.3	84.3
XLM-100	79.9	91.6	93.7	91.8	42.5	85.6	69.8	63.0	96.8	84.2
XLM-17	80.2	91.9	93.7	92.0	44.8	86.3	70.6	61.8	96.3	84.5
HerBERT	80.5	92.7	92.5	91.9	50.3	89.2	76.3	52.1	95.3	84.5
XLM-R base	81.5	92.1	94.1	93.3	51.0	89.5	74.7	55.8	98.2	85.2
Polbert	81.7	93.6	93.4	93.8	52.7	87.4	71.1	59.1	98.6	85.2
Our model	85.3	93.9	94.2	94.0	66.7	90.6	76.3	65.9	98.8	87.8
Large models										
XLM-R large	87.5	94.1	94.4	94.7	70.6	92.4	81.0	72.8	98.9	88.4
Our model	87.8	94.5	93.3	94.9	71.1	92.8	82.4	73.4	98.8	88.8

dom Fields (CRF) inference layer, and multiple Viterbi-decoding steps to handle nested entities. In our experiments, we used the same hyperparameters as the authors.

3.3 Results and discussion

In this section, we demonstrate the results of evaluating our language models on downstream tasks. We repeated the fine-tuning of the models for each task five times. The scores reported are the median values of those five runs. Table 1 demonstrates the evaluation results on the KLEJ benchmark, in comparison with other available Polish and multilingual transformer-based models. The results of other approaches are taken from the KLEJ leaderboard. We split the table into two sections, comparing the BERT-base and BERT-large architectures separately. We can observe that there is a wider selection of base models, and most of them are multilingual, such as the original multilingual BERT (mBERT) [16], SlavicBERT [4], XLM [7], and XLM-R [8]. The only models pre-trained specifically for Polish language are HerBERT [39] and Polbert. Among the base models, our approach outperforms others by a significant margin. In the case of large models, only the XLM-RoBERTa (XLM-R) pre-trained model has been available until now. XLM-RoBERTa is a recently published multilingual transformer trained on 2.5TB of data in 100 languages. It has been shown to be highly competitive against monolingual models. A direct comparison with our Polish language model demonstrates a consistent advantage of our model - it has achieved better results in seven of the nine tasks included in the KLEJ benchmark.

Table 2 shows a more detailed breakdown of the evaluation results, and includes all the tasks from the KLEJ benchmark, and four additional tasks: SICK-R, SICK-R, 8TAGS, and PolEval-NER 2018. For each task, we define the task type (classification, regression, or sequence tagging), the metric used for evaluation, the previous state-of-the-art, and our results including the absolute difference to the SOTA. The competition between XLM-R and our large model dominates the results, since both models have led to significant improvements in linguistic tasks for Polish language. In some cases, the improvement over previous approaches is greater than 10%. For example, the CDB task was a part of the PolEval 2019 competition, in which the winning solution by Czapla et al. [10] achieved an F1-score of 58.6. Both our model and the XLM-R large model outperform that by at least twelve points, achieving an F1-score of over 70. The comparison for the named entity recognition task is also interesting. The previous state-of-the-art solution by Dadas [11] is a model that combined neural architecture with external knowledge sources, such as entity lexicons or a specialized entity linking module based on data from Wikipedia. Our language model managed to outperform this method by 3.8 points without using any structured external knowledge. In summary, our model has demonstrated an improvement over existing methods in eleven of the thirteen tasks.

4 Conclusions

We have presented two transformer-based language models for Polish, pre-trained using a combination of publicly available text corpora and a large collection of methodically pre-processed web data. We have shown the effectiveness

Table 2: Detailed results for Polish language downstream tasks. In some cases, we used the datasets and task definitions from the KLEJ benchmark, which are different from the original tasks they were based on (they have been reformulated or otherwise modified by the benchmark authors). We denote such tasks with (KLEJ) to emphasize that the evaluation was performed on the KLEJ version of the data. The abbreviated task types are: C - classification, R - regression, and ST - sequence tagging.

Task		Metric	Previous SOTA		Base model	Large model
Multi-class classification						
NKJP (KLEJ)	C	Accuracy	XLM-R large [8]	94.1	93.9 (−0.2)	94.5 (+0.4)
8TAGS	C	Accuracy	ELMo [12]	71.4	77.2 (+5.8)	80.8 (+9.4)
Binary classification						
CBD	C	F1-score	XLM-R large [8]	70.6	66.7 (−2.9)	71.1 (+0.5)
DYK (KLEJ)	C	F1-score	XLM-R large [8]	72.8	65.9 (−6.9)	73.4 (+0.6)
PSC (KLEJ)	C	F1-score	XLM-R large [8]	98.9	98.8 (−0.1)	98.8 (−0.1)
Sentiment analysis						
PolEmo2.0-IN	C	Accuracy	XLM-R large [8]	92.4	90.6 (−1.8)	92.8 (+0.4)
PolEmo2.0-OUT	C	Accuracy	XLM-R large [8]	81.0	76.3 (−4.7)	82.4 (+1.4)
Allegro Reviews	R	1-wMAE	XLM-R large [8]	88.4	87.8 (−1.0)	88.8 (+0.4)
Textual entailment						
CDSC-E	C	Accuracy	XLM-R large [8]	94.4	94.2 (−0.2)	93.3 (−1.1)
SICK-E	C	Accuracy	LASER [12]	82.2	86.1 (+3.9)	87.7 (+5.5)
Semantic relatedness						
CDSC-R	R	Spearman	XLM-R large [8]	94.7	94.0 (−0.7)	94.9 (+0.2)
SICK-R	R	Spearman	USE [12]	75.8	82.3 (+6.5)	85.6 (+9.8)
Named entity recognition						
Poleval-NER 2018	ST	F1-score	Dadas [11]	86.2	87.9 (+1.7)	90.0 (+3.8)

of our models by comparing them with other transformer-based approaches and recent state-of-the-art approaches. We conducted a comprehensive evaluation on a wide set of Polish linguistic tasks, including binary and multi-class classification, regression, and sequence labeling. In our experiments, the larger model performed better than other methods in eleven of the thirteen cases. To accelerate research on NLP for Polish language, we have released the pre-trained models publicly.

References

- [1] Alan Akbik, Duncan Blythe, and Roland Vollgraf. Contextual string embeddings for sequence labeling. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1638–1649, 2018.
- [2] Estera Małek Aleksander Wawer. Results of the PolEval 2018 Shared Task 2: Named Entity Recognition. *Proceedings of the PolEval 2018 Workshop*, pages 53–62, 2018.
- [3] Wissam Antoun, Fady Baly, and Hazem Hajj. Arabert: Transformer-based model for arabic language understanding. *arXiv preprint arXiv:2003.00104*, 2020.
- [4] Mikhail Arhipov, Maria Trofimova, Yuri Kuratov, and Alexey Sorokin. Tuning multilingual transformers for language-specific named entity recognition. In *Proceedings of the 7th Workshop on Balto-Slavic Natural Language Processing*, pages 89–93, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-3712. URL <https://www.aclweb.org/anthology/W19-3712>.
- [5] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146, 2017.

- [6] José Cañete, Gabriel Chaperon, Rodrigo Fuentes, and Jorge Pérez. Spanish pre-trained bert model and evaluation data. In *Practical ML for Developing Countries Workshop @ ICLR 2020*, 2020.
- [7] Alexis Conneau and Guillaume Lample. Cross-lingual language model pretraining. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 7059–7069. Curran Associates, Inc., 2019. URL <http://papers.nips.cc/paper/8928-cross-lingual-language-model-pretraining.pdf>.
- [8] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzman, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*, 2019.
- [9] Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, Ziqing Yang, Shijin Wang, and Guoping Hu. Pre-training with whole word masking for chinese bert. *arXiv preprint arXiv:1906.08101*, 2019.
- [10] Piotr Czapla, Sylvain Gugger, Jeremy Howard, and Marcin Kardas. Universal language model fine-tuning for polish hate speech detection. *Proceedings of the PolEval 2019 Workshop*, page 149, 2019.
- [11] Slawomir Dadas. Combining neural and knowledge-based approaches to named entity recognition in polish. *International Conference on Artificial Intelligence and Soft Computing*, pages 39–50, 2019.
- [12] Slawomir Dadas, Michał Perelkiewicz, and Rafał Poświata. Evaluation of Sentence Representations in Polish. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pages 1674–1680, Marseille, France, May 2020. European Language Resources Association. ISBN 979-10-95546-34-4. URL <https://www.aclweb.org/anthology/2020.lrec-1.207>.
- [13] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc Le, and Ruslan Salakhutdinov. Transformer-XL: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2978–2988, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1285. URL <https://www.aclweb.org/anthology/P19-1285>.
- [14] Wietse de Vries, Andreas van Cranenburgh, Arianna Bisazza, Tommaso Caselli, Gertjan van Noord, and Malvina Nissim. Bertje: A dutch bert model. *arXiv preprint arXiv:1912.09582*, 2019.
- [15] Pieter Delobelle, Thomas Winters, and Bettina Berendt. Robbert: a dutch roberta-based language model. *arXiv preprint arXiv:2001.06286*, 2020.
- [16] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. 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 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://www.aclweb.org/anthology/N19-1423>.
- [17] Kenneth Heafield. Kenlm: Faster and smaller language model queries. *Proceedings of the Sixth Workshop on Statistical Machine Translation*, pages 187–197, 2011.
- [18] Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 328–339, 2018.
- [19] Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. What does BERT learn about the structure of language? In *ACL 2019 - 57th Annual Meeting of the Association for Computational Linguistics*, Florence, Italy, July 2019. URL <https://hal.inria.fr/hal-02131630>.
- [20] Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *International Conference on Learning Representations*, 2020.
- [21] Jan Kocoń, Piotr Miłkowski, and Monika Zaśko-Zielińska. Multi-level sentiment analysis of PolEmo 2.0: Extended corpus of multi-domain consumer reviews. In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 980–991, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/K19-1092. URL <https://www.aclweb.org/anthology/K19-1092>.

- [22] Olga Kovaleva, Alexey Romanov, Anna Rogers, and Anna Rumshisky. Revealing the dark secrets of BERT. 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 4365–4374, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1445. URL <https://www.aclweb.org/anthology/D19-1445>.
- [23] Taku Kudo and John Richardson. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium, November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-2012. URL <https://www.aclweb.org/anthology/D18-2012>.
- [24] Yuri Kuratov and Mikhail Arkhipov. Adaptation of deep bidirectional multilingual transformers for russian language. *arXiv preprint arXiv:1905.07213*, 2019.
- [25] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations. In *International Conference on Learning Representations*, 2020.
- [26] Hang Le, Loïc Vial, Jibril Frej, Vincent Segonne, Maximin Coavoux, Benjamin Lecouteux, Alexandre Allauzen, Benoît Crabbé, Laurent Besacier, and Didier Schwab. Flaubert: Unsupervised language model pre-training for french. *arXiv preprint arXiv:1912.05372*, 2019.
- [27] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [28] Michał Marcińczuk, Marcin Ptak, Adam Radziszewski, and Maciej Piasecki. Open Dataset for Development of Polish Question Answering Systems. In Z. Vetulani and H. Uszkoreit, editors, *Proceedings of Human Language Technologies as a Challenge for Computer Science and Linguistics’13*, pages 479–483, Poznań, 2013. Fundacja UAM.
- [29] Marco Marelli, Stefano Menini, Marco Baroni, Luisa Bentivogli, Raffaella Bernardi, and Roberto Zamparelli. A SICK cure for the evaluation of compositional distributional semantic models. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, pages 216–223, Reykjavik, Iceland, May 2014. European Language Resources Association (ELRA). URL http://www.lrec-conf.org/proceedings/lrec2014/pdf/363_Paper.pdf.
- [30] Louis Martin, Benjamin Muller, Pedro Javier Ortiz Suárez, Yoann Dupont, Laurent Romary, Éric Villemonte de la Clergerie, Djamel Seddah, and Benoît Sagot. Camembert: a tasty french language model. *arXiv preprint arXiv:1911.03894*, 2019.
- [31] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.
- [32] Dat Quoc Nguyen and Anh Tuan Nguyen. Phobert: Pre-trained language models for vietnamese. *arXiv preprint arXiv:2003.00744*, 2020.
- [33] Maciej Ogrodniczuk and Mateusz Kopeć. The polish summaries corpus. In Nicoletta Calzolari (Conference Chair), Khalid Choukri, Thierry Declerck, Hrafn Loftsson, Bente Maegaard, Joseph Mariani, Asuncion Moreno, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, Reykjavik, Iceland, may 2014. European Language Resources Association (ELRA). ISBN 978-2-9517408-8-4.
- [34] Tony Peng. The staggering cost of training sota ai models, Jun 2019. URL <https://syncedreview.com/2019/06/27/the-staggering-cost-of-training-sota-ai-models/>.
- [35] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.

- [36] Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, volume 1, pages 2227–2237, 2018.
- [37] Adam Przepiórkowski, Mirosław Banko, Rafał L Górski, and Barbara Lewandowska-Tomaszczyk. Narodowy Korpus Języka Polskiego [Eng.: National Corpus of Polish]. *Wydawnictwo Naukowe PWN, Warsaw*, 2012.
- [38] Michal Ptaszynski, Agata Pieciukiewicz, and Paweł Dybała. Results of the PolEval 2019 Shared Task 6: First Dataset and Open Shared Task for Automatic Cyberbullying Detection in Polish Twitter. *Proceedings of the PolEval 2019 Workshop*, page 89, 2019.
- [39] Piotr Rybak, Robert Mroczkowski, Janusz Tracz, and Ireneusz Gawlik. KLEJ: Comprehensive Benchmark for Polish Language Understanding. *arXiv preprint arXiv:2005.00630*, 2020.
- [40] Takashi Shibuya and Eduard Hovy. Nested named entity recognition via second-best sequence learning and decoding. *arXiv preprint arXiv:1909.02250*, 2019.
- [41] Fabio Souza, Rodrigo Nogueira, and Roberto Lotufo. Portuguese named entity recognition using bert-crf. *arXiv preprint arXiv:1909.10649*, 2019. URL <http://arxiv.org/abs/1909.10649>.
- [42] Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Hao Tian, Hua Wu, and Haifeng Wang. Ernie 2.0: A continual pre-training framework for language understanding. *arXiv preprint arXiv:1907.12412*, 2019.
- [43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [44] Antti Virtanen, Jenna Kanerva, Rami Ilo, Jouni Luoma, Juhani Luotolahti, Tapio Salakoski, Filip Ginter, and Sampo Pyysalo. Multilingual is not enough: Bert for finnish. *arXiv preprint arXiv:1912.07076*, 2019.
- [45] Alina Wróblewska and Katarzyna Krasnowska-Kieraś. Polish evaluation dataset for compositional distributional semantics models. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 784–792, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1073. URL <https://www.aclweb.org/anthology/P17-1073>.
- [46] Liang Xu, Xuanwei Zhang, and Qianqian Dong. CLUECorpus2020: A large-scale chinese corpus for pre-training language model. *arXiv preprint arXiv:2003.01355*, 2020.
- [47] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. In *Advances in neural information processing systems*, pages 5754–5764, 2019.

A Non-English transformer-based language models

Project	Languages	Paper	Project URL
Arabic-BERT	Arabic	-	github.com/alisafaya/Arabic-BERT
AraBERT	Arabic	[3]	github.com/aub-mind/arabert
ClueCorpus2020	Chinese	[46]	github.com/CLUEbenchmark/CLUECorpus2020
Chinese BERT	Chinese	[9]	github.com/ymcui/Chinese-BERT-wwm
Google BERT	Chinese	[16]	github.com/google-research/bert
ERNIE 2.0	Chinese	[42]	github.com/PaddlePaddle/ERNIE
BERTje	Dutch	[14]	github.com/wietsedv/bertje
RoBBERT	Dutch	[15]	ipieter.github.io/blog/robbert
Finnish BERT	Finnish	[44]	github.com/TurkuNLP/FinBERT
CamemBERT	French	[30]	camembert-model.fr
FlauBERT	French	[26]	github.com/getalp/Flaubert
German BERT	German	-	deepset.ai/german-bert
GreekBERT	Greek	-	github.com/nlpaueb/greek-bert
UmBERTo	Italian	-	github.com/musixmatchresearch/umberto
GilBERTo	Italian	-	github.com/idb-ita/GilBERTo
Japanese BERT	Japanese	-	github.com/yoheikikuta/bert-japanese
Japanese BERT	Japanese	-	github.com/cl-tohoku/bert-japanese
KoBERT	Korean	-	github.com/SKTBrian/KoBERT
Malaya	Malaysian	-	github.com/huseinzol05/Malaya/
Nordic BERT	Nordic (4)	-	github.com/botxo/nordic_bert
PolBERT	Polish	-	github.com/kldarek/polbert
HerBERT	Polish	[39]	klejbenchmark.com
Portuguese BERT	Portuguese	[41]	github.com/neuralmind-ai/portuguese-bert
RuBERT	Russian	[24]	github.com/deepmipt/DeepPavlov
SlavicBERT	Slavic (4)	[4]	github.com/deepmipt/Slavic-BERT-NER
BETO	Spanish	[6]	github.com/dccuchile/beto
Swedish BERT	Swedish	-	github.com/Kungbib/swedish-bert-models
BERTsson	Swedish	-	huggingface.co/jannesg/bertsson
BERTurk	Turkish	-	github.com/stefan-it/turkish-bert
PhoBERT	Vietnamese	[32]	github.com/VinAIResearch/PhoBERT