# Fine-tuning Transformer-based Encoder for Turkish Language Understanding Tasks

SAVAŞ YILDIRIM, Istanbul Bilgi University, Faculty of Engineering and Natural Sciences

Deep learning-based and lately Transformer-based language models have been dominating the studies of natural language processing in the last years. Thanks to their accurate and fast fine-tuning characteristics, they have outperformed traditional machine learning-based approaches and achieved state-of-the-art results for many challenging natural language understanding (NLU) problems. Recent studies showed that the Transformer-based models such as BERT, which is Bidirectional Encoder Representations from Transformers, have reached impressive achievements on many tasks. Moreover, thanks to their transfer learning capacity, these architectures allow us to transfer pre-built models and fine-tune them to specific NLU tasks such as question answering. In this study, we provide a Transformer-based model and a baseline benchmark for the Turkish Language. We successfully fine-tuned a Turkish BERT model, namely BERTurk that is trained with base settings, to many downstream tasks and evaluated with a the Turkish Benchmark dataset. We showed that our studies significantly outperformed other existing baseline approaches for Named-Entity Recognition, Sentiment Analysis, Question Answering and Text Classification in Turkish Language. We publicly released these four fine-tuned models and resources in reproducibility and with the view of supporting other Turkish researchers and applications.

## 1 INTRODUCTION

The advent of deep learning and transfer learning in the field of natural language processing (NLP) has achieved significant improvements for many tasks. Sequence transduction models such as RNN and LSTM are based on complex recurrent neural networks. They successfully connected the encoder and decoder for sequence to sequence problem [31]. However one of the computational bottlenecks of RNN-based models is processing very long sequences and they have difficulty in keeping long-distance dependencies.

The Transformers [23] overcome this limitation and simplify encoder-decoder network architecture by employing attention mechanisms, dropping recurrence or convolutions entirely. The studies showed that the Transformer could be trained significantly faster than other standard RNN architectures since that RNN models are based on more complex recurrent layers. The initial experiments on some machine translation tasks indicated that the Transformer-based models are superior in quality. They are more parallelizable and further requiring significantly less time to train. Finally, the Transformer-based models became a new paradigm in NLP and led to many

Author's address: Savaş Yıldırım, savas.yildirim@bilgi.edu.tr, Istanbul Bilgi University, Faculty of Engineering and Natural Sciences, Eski Silahtarağa Elektrik Santralı Kazım Karabekir Cad. No: 2/13, Istanbul, Turkey, 34060.

successful derivatives that provides state-of-the-art general-purpose architectures tested on General Language Understanding Evaluation (GLUE) Benchmark tasks [24]; BERT [6], GPT-2 [18], RoBERTa [13], XLM [11], XLNet [26], T5 [19] and so forth.

Main advantage of the transformer models is to provide fast and cheap fine-tuning in terms of space and time complexity by transferring the base language models trained in the pre-training phase. While the pre-training phase demands huge amounts of data, the fine-tuning phase requires relatively small amounts of data depending on downstream tasks. The pre-training phase trains the language model so that it can build internal representations such as contextual word embeddings or sentence encoding [17], which allows there after a model to be reused for different downstream tasks even with little amount of labeled data. The contextualized word embeddings [17], addressed the well known word sense problem and represent the word within its context contrary to other traditional word embeddings like word2vec [14], GloVe [16], fastText[1]. Another advantage of the attention mechanism allows the model to better resolve long-term dependencies and coreference relation between entities [8]. This property is important especially for machine translation and question answering problems.

In this study we provide a benchmark for Natural Language Understanding (NLU) downstreams tasks for the Turkish Language. We fine-tuned a pre-built BERT language model, BERTurk [20], to four downstream tasks; *Named-Entity Recognition, Sentiment Analysis, Question Answering and Text Classification*. Compared to the baseline approaches, we have got successful results in the NLU tasks. In order to reproduce the result and support the community, we publicly released these four fine-tuned models. This study is considered as the first successful attempt to fine-tune the BERT model for these four tasks in the Turkish Language.

## 2 RELATED WORKS FOR LANGUAGE MODELS

The traditional language models have been based on n-grams and they were simply a probability distribution over word sequences with an objective function of predicting the last *n-th word* from previous *(n-1) words*. One drawback of n-gram language model is that as the model is being trained with a very large corpus, the number of unique words increases and therefore sparse data problem arises [31].

Modern neural networks sorted the dimensionality problem out by using continuous and dense representation for the words, called *word embeddings*. The early neural language representation employed encoders to produce static word embeddings. The most popular models, such as *word2vec [14], glove [16] and fastText [1]* transformed unsupervised textual data into a supervised one by either predicting the target word using context or predicting contextual words by the target word based on a sliding window of n-gram context. The final output turns out to be context-independent word embeddings which ignores word sense.

The main problem of these early models is that they never use the sentence level representation. Second problem is that the senses of the words as one of the oldest problems in computational linguistics are ignored and a single fixed representation is assigned to each word or token. Further studies used a neural network component which encodes the complete sentence-level information and *contextual word embeddings* which finally alleviate word sense problems. In early days, the essential architectures such as RNN, CNN, LSTM have been widely used as encoder and decoder in a sequence to sequence problems.

*ELMo (Embeddings from Language Models)* representations [17], which is based on LSTM, successfully built *deep contextualized word representations* and solved the sense and sentence encoding problem. The ELMo representations differ from traditional word embeddings in that each word representation is a function of the entire input sentence.

Later on, the Transformers [23] simplified sequence-to-sequence models by removing recurrency and only keeping attention. The Transformer-based architecture such as GPT [18] and BERT [6] successfully captured the contextualized sentence-level language representations and word embeddings in a modern way. These sentence encoders differ from sliding window approaches in that they read a complete sentence instead of any fixed sized word sequences defined by a sliding window. While Elmo used bidirectional LM to read the entire sentence independently in both directions and concatenates trained forward and backward LSTMs, BERT adopted Transformer to read and process the entire input at once in a simultaneous bidirectional way.

BERT algorithm is able to read the input in a bidirectional way rather than unidirectional. The BERT model trains two different objectives simultaneously; *Masked Language Model (MLM)* and *Next Sentence Prediction (NSP)*. MLM is based on optimizing the prediction of the masked words in sentences that are randomly chosen. MLM randomly replaces some masked tokens with a special symbol, such as ***[MASK]***. Then, the model training objective is to predict the masked symbol based on the non-masked context. The major difference between Transformer-based BERT and LSMT-based ELMo is that the predictions in BERT are based on the entire context rather than only one direction. The BERT is however criticized that the masks could be never seen at fine-tuning, which faces a mismatch between pre-training and fine-tuning phases. NSP is based on a model that detects whether two sentences follow each other or not. Some other transformer models such as RoBERTa [13] or Albert [12] employ different approaches for the objective functions to improve the results.

Although both RNN and Transformer architectures gained spectacular achievements, the main limitation is said to be understanding long-term dependencies. Transformer-XL [3] put together the advantages of these two architectures. It employs the self-attention modules on each contextual segment of input and a RNN mechanism to learn dependencies between successive segments. Based on the inspirations of impressive RNN and Transformer-based models like ELMo [17] and BERT [6], recently a wide range of LM derivatives have been proposed. The variants such as *Electra [2], Albert [12], T5 [19], RoBERTa [13]* are gaining successful achievements on GLUE tasks [24] which is a set of tools for evaluating the performance of models across a diverse collection of common NLU tasks.

## 3 THE TRANSFORMER-BASED MODEL FOR THE TURKISH LANGUAGE

The advantage of the Transformers is that the same architecture could be used in both pre-training and fine-tuning except output layers. The pre-trained model parameters are used to initialize models for different downstream tasks. Pre-training and fine-tuning made it possible to apply transfer learning in NLP and have become a new paradigm. In this section, we describe how to fine-tune BERT-based models for a variety of problems in Turkish; *Question Answering, Sentiment Analysis, Named-Entity Recognition and Text classification*. We also discuss how to evaluate them with a set of benchmark datasets.

### 3.1 Fine-tuning BERT model for Turkish

We fine-tuned four NLU downstream tasks based on *BERTurk* [20] and uploaded to the repository. BERTurk [20] has been successfully pre-trained as a Transformer model with a huge Turkish corpus and shared with the NLP community. There are two alternatives for BERT training; *base and large*. While base architecture has *(L=12, H=768, A=12, Total Parameters=110M)*, large one has *(L=24, H=1024, A=16, Total Parameters=340M)* where L is number of layers, and H is hidden size, A is number of self attention heads. The model can be used *cased* or *uncased* as a way of formatting input text. For the Turkish language, we have currently only base pre-trained LM of BERT with cased or uncased.

The last version of Turkish BERT model, which we used in this study, has been trained on a huge training corpus with file size of 35GB and 45B tokens. The compiled corpus consists of online resources such as *OSCAR Corpus* [1], *OPUS corpus* [2], *Wikipedia dump* and *Turkish Tree-bank* [3]. The base model has been trained within a cloud environment provided by *Google TensorFlow Research Cloud (TFRC)*. Uncased and cased models with *128K* and *32K* Vocabulary have been successfully trained. In order to load or fine-tune the model, we used *python* and its *transformer* library [4], which is built for the NLP community to share, load and extend their models.

The Transformers can learn syntactic segmentation at the lower levels such as *subwords* and then semantic knowledge at the higher levels such as *coreference resolution* as they work on specific language. BERT model applies *WordPiece* tokenization procedure to each token to understand syntactic knowledge. The implementation is simply based on the one from *tensor2tensor*, WordPiece embeddings [25]. The *wordpiece* model is generated using a data-driven approach to maximize the language-model likelihood of corpus. When a training corpus is provided, the optimization problem is to select limited wordpieces. Turkish is an agglutinative language and the root word receives many suffix so that incredibly large word could be generated as follows (wikipedia.org)

```
"muvaffakiyetsizleştiricileştiriveremeyebileceklerimizdenmişsiniz"
```

This word is one of the longest words and can be extended more. It is derived from the noun word *"muvaffakiyet"* (success). It means *"As though you are from those whom we may not be able to easily make into a maker of unsuccessful ones"*. The response of wordpiece tokenization procedure to the word is truly successful as follows

```
['muvaffak', '##iyet', '##siz', '##leş', '##tir', '##ici',
 '##leş', '##tir', '##iver', '##emeye', '##bilecekleri', '##mi', '##z',
 '##den', '##miş', '##siniz', '##cesine', '##siniz']
```

It is now transformed into 18 chunks and represented by (18+2, 32000) tensors in the transformer systems where +2 is needed for [CLS] and [SEP] tokens. Such capacity of the syntactic breakdown helps many models to easily fine-tune downstream tasks. We fine-tuned four models that are now easily accessible via *Hugging-Face platform* [5].

### 3.2 Sentiment Analysis

The fine-tuned model has been already shared [6]. The dataset for sentiment analysis is obtained by two studies [5], [7]. Erkin Demirtas and Mykola Pechenizkiy [5] gathered movie and product reviews from popular e-commerce and social web sites. The set is gathered from a popular Turkish cinema web page [7] with over 5K positive and 5K negative comments in the Turkish language. Reviews are marked in scale from 0 to 5 by the web users. The study considered a review sentiment positive if the rating is equal to or bigger than 4, and negative if it is less or equal to 2 as the study does. The products review dataset is taken from an online e-commerce web page. This benchmark dataset consists of reviews regarding four product categories. The categories of

---

[1] https://traces1.inria.fr/oscar/

[2] http://opus.nlpl.eu/

[3] https://ii.metu.edu.tr/metu-corpora-research-group

[4] https://huggingface.co/transformers/

[5] https://huggingface.co/savasy

[6] https://huggingface.co/savasy/bert-base-turkish-sentiment-cased

[7] www.beyazperde.com

Table 1. Dataset for Sentiment Analysis

| Sentiment Dataset | # of total comments | Twitter | E-commerce (with Movie) | #Pos | #Neg |
|---|---|---|---|---|---|
| Training Set | 40830 | 22000 | 18830 | 20340 | 20490 |
| Development Set | 9000 | 4900 | 4100 | 4445 | 4556 |
| Test Set | 9000 | 4900 | 4100 | 4391 | 4609 |

Table 2. Sentiment Analysis Performance Evaluation

| Model | Dataset | Results |
|---|---|---|
| Yildirim (2020) | e-commerce | 89.26 |
| Demirtas (2013) | e-commerce | 85.1 |
| Sert (2017) [7] | Twitter | 80.5 |
| *Our BERT model** | e-commerce | **93.12** |
| *Our BERT model** | Twitter | **96.69** |
| *Our BERT model** | e-commerce + Twitter | **94.77** |

e-commerce dataset are *dvd, books, electronics and kitchen*. Likewise, reviews are marked in the range from 1 to 5. Each category has equally 700 positive and 700 negative reviews, where average rating of negative reviews is 2.27 and of positive reviews is 4.5. The dataset has been evaluated in another study [30] as well. Finally we have organized all the datasets for training, development and test set as described in Table 1

The sentiment model was fine-tuned with max-seq-length of 128, learning rate of 2e-5, number of epoch optimized is 3 [8]. The performance of the model is reported as shown in Table 2. We listed other studies performance that used the same e-commerce and Twitter datasets. We only report their best performance for a proper comparison. We got *93.12%* and *96.69%* accuracy for e-commerce and Twitter dataset respectively. When we merged two datasets for training and test phase, we got *94.77%* accuracy .

## 3.3 Named-Entity Recognition

The fine-tuned model can be found under the link [9] Named-Entity Recognition aims to capture named entities in a text, such as locations, person. Our NER datasets follow *BIO (Beginning, Inside, Outside)* data format which is a common NER tagging format for tokens. The *B* refers to the first word of an entity and the *I* corresponds to the remaining words of the same entity. Tokens tagged with *O* means are not part of an entity. B and I tags are followed by an entity category such as *B-PER, I-PER, B-ORG*. NER problem is considered a multi-class token classification task. There are a couple of named-entity data formats: *Enamex, Timex, Numex. Enamex*, which is used in our system, has three entity types: Person, location, organization. NUMEX format is for numerical entities such as money and percent and TIMEX format is used for temporal entities such as date and time. Two different datasets annotated in Enamex format have been used for an evaluation in our fine-tuning phase as shown in Table 3. The first one is Turkish NER dataset obtained from the study *WikiANN* dataset [15] which includes 282 different language dataset compiled from Wikipedia. Second one which is also annotated with POS, LOC, and ORG is shared by NLP community [20, 27]

---

[8]https://huggingface.co/savasy/bert-base-turkish-sentiment-cased
[9]https://huggingface.co/savasy/bert-base-turkish-ner-cased

Table 3. NER Dataset

| Model | Dataset | # of token | # PER | # LOC | # ORG |
|---|---|---|---|---|---|
| BERT Model | WikiANN | 340K | 17.2K | 19.1K | 15.8K |
| BERT Model | Second Dataset | 385K | 16.6K | 10.8K | 10.1K |

Table 4. Performance Evaluation of NER Models

| Model | Dataset | F1 Score |
|---|---|---|
| Pan (2017) | WikiANN | **96.9** |
| Kucuk (2014) | Twitter data | 72.61 |
| Yeniterzi (2018) | Community Set | 88.94 |
| Demir (2014) | by [22] | 91.85 |
| Seker (2012) | by [22] | 91.94 |
| Seker (2017) | by [22] | 92.00 |
| Tur (2003) | General News | 92.73 |
| Our BERT Model* | WikiANN | **95.0** |
| Our BERT Model* | Community Set | **95.2** |

The initial studies of Turkish NER have started in 1990 [28]. *Hidden-markov models, Conditional Random Fields, Bayesian Models and Neural Networks* are among the approaches applied to the NER solution. We listed NER performances of some Turkish studies [4, 9, 15, 21, 22, 27, 32] that especially uses Enamex datasets as shown in Table 4. Bert-based fine-tuning model gained comparable results of 95.0 F1 and 95.2 F1 score for two datasets as reported in the table. The study [15] still has the highest score of *96.9 F1* by leveraging cross-lingual language models. They developed a cross-lingual name tagging and linking framework for 282 different languages that exist in Wikipedia. The framework could identify name mentions and assign a coarse-grained type to the mentions. It applies Bi-directional Long Short-Term Memory and Conditional Random Fields (CRFs) network as the underlying learning model for the name tagger for each language. It also links the captured names to an English Knowledge Base.

## 3.4 Question Answering
The fine-tuned model can be found under the link[10]

Turkish question answering dataset, namely *TQuAD* [11], has been prepared in the form of *Stanford Question Answering Dataset (SQuAD)* which is a well-known English QA dataset. The SQuAD includes questions on a set of Wikipedia articles, where the answer to a question is a span of a text from the corresponding reading passage. While SQuAD 1.0 format does not cover features that the question might be *unanswerable*, this feature was later added on the *SQuAD 2.0* version. The *TQuAD dataset* currently follows SQuAD 1.0. The examples have been manually prepared on *Turkish & Islamic Science History* within the scope of *Teknofest 2018 Artificial Intelligence competition*, and it has been shared with the NLP community. Training and evaluation set have *681 and 72 titles, 2232 and 275 Paragraphs, and 8308 and 892 Q&As* respectively. The model has been fine-tuned with the default BERT model settings such as learning rate of 0.00003, maximum sequence

[10]https://huggingface.co/savasy/bert-base-turkish-squad
[11]https://github.com/TQuad/turkish-nlp-qa-dataset

Table 5. Text Classification Performance

| Model | Dataset | Size | *Performance* |
|---|---|---|---|
| Yildirim (2014) | TTC 4900 7 Classes | 4900 | 90.0 |
| Kilinc (2017) | TTC 3600 6 Classes | 3600 | 91.3 |
| Our Fine-tuned BERT* | TTC 4900 7 Classes | 4900 | **93.4** |
| Our Fine-tuned BERT* | TTC 3600 6 Classes | 3600 | **92.2** |

length of 384. We observed that model successfully fine-tuned the base pre-trained model with the exact match score of *62.55* and F1 score of **80.48**. For English Language we can see above *95.0 F1* score and we have enough room for enhancing the Question Answering model by increasing data size and hyper-parameter optimization. Due to the fact that there is no study using SQuAD-like dataset in Turkish, we are not able to compare our model to any Turkish study.

## 3.5 Text Classification

The fine-tuned model can be found under the link [12]. Text classification downstream task has been fine-tuned with two datasets *TTC-4900* [29] and *TTC-3600* [10]. TTC-4900 is originally prepared by *a Turkish NLP Group* [13] and shared with community via repository [14]. The data has been already used by many studies [29]. The dataset has seven coarse-grained categories: *world, economy, culture, health, politics, sport, and technology*. The dataset equally consists of 700 articles under each category. *TTC-3600* dataset has been shared by the study [10] with a link [15]. It has six categories; *economy, culture, health, politics, sports, technology*. Each category has 600 examples. The BERT-based model has been fine tuned with these two datasets in the same settings with the original BERT model and reported as shown in Table 5. As the table indicates, our fine-tuned models surpassed the traditional approaches employed in the other studies [10, 29] in terms of F1 score by achieving **93.4** and **92.2**.

## 4 FINAL REMARKS

We evaluated four different downstream tasks for the Turkish language. We gathered several datasets used and shared by the community. Other than Question Answering dataset, the remaining three datasets have been already used and evaluated by other studies. We compared our fine-tuned models with those studies who used the same datasets as shown in Table 2, Table 4, and Table 5, which suggests that the BERT-based models outperformed many baseline approaches implemented by other studies.

We observed some well-known facts that when we transfer the model to another domain, we obtain less performance. For instance, we measure the sentiment prediction scores ranging from 95.0 F1 to 80.0 F1. But when we use the same data from the same domain, the performance improves. As we keep fine-tuning the process fed by that particular domain, the model gets again over 95.0 % F1 score. The Named-Entity Recognition model follows ENAMEX format where three coarse-grained entity types are considered: PER, LOC, ORG. We want to train NER models for the problems of NUMEX and TIMEX format. Our preliminary studies indicated that we can get high performances around 93 F1 with these formats, which is not reported here. The study [15] has an impressive F1 score of 96.9 and better than our models.

---

[12] savasy/bert-turkish-text-classification

[13] http://www.kemik.yildiz.edu.tr

[14] http://www.kaggle.com/savasy/ttc4900

[15] https://github.com/denopas/TTC-3600

Question Answering model cannot cover unanswerable questions, which means no-answer option is not available and always returns a most likely answer, which is a standard for SQuAD 1.0. We will extend our model to SQuAD 2.0 by adding a "no-answer" option. When comparing our QA study to English counterpart, one can say that both the size of Turkish dataset and the performance of the model needs to be improved. In the last decades, Text Classification problem has been successfully solved by traditional models, TF-IDF, or word embeddings [10, 29] in Turkish. We showed that the BERT-based model has outperformed the traditional approaches for both seven-class and six-class classification problems. In the future we will fine-tune a model to cover more fine-grained classes.

There exist many other NLU problems to be addressed in the Turkish Language as one of the less studied languages. We also plan to leverage other transformer models such as RoBERTa or ELECTRA to improve the performance of NLU tasks and solve other types of NLU problems such as anaphora resolution, summarization, relation extraction in Turkish.

## 5   CONCLUSION

The recent studies showed that transformer-based models got state-of-the-art results in many downstream NLP tasks. In this work, we have leveraged them for four different downstream NLU tasks: Named-Entity Recognition, Sentiment Analysis, Question Answering and Text Classification. The experiments showed that the fine-tuned models significantly outperformed other existing baseline approaches in the Turkish Language as reported in the study. We publicly released these four fine-tuned models in reproducibility and in supporting other researchers. This study is counted to be the first successful attempt to apply the transformer-based model to the Turkish language by fine-tuning the BERT-based model. In the future, we want to leverage other types of language models and cover more challenging down-streams tasks such as coreference resolution or text summarization for the Turkish language.

## REFERENCES

[1] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2016. Enriching Word Vectors with Subword Information. *arXiv preprint arXiv:1607.04606* (2016).

[2] Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. 2020. Electra: Pre-training text encoders as discriminators rather than generators. *arXiv preprint arXiv:2003.10555* (2020).

[3] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G. Carbonell, Quoc V. Le, and Ruslan Salakhutdinov. 2019. Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context. *CoRR* abs/1901.02860 (2019). arXiv:1901.02860 http://arxiv.org/abs/1901.02860

[4] H. Demir and A. Özgür. 2014. Improving Named Entity Recognition for Morphologically Rich Languages Using Word Embeddings. In *2014 13th International Conference on Machine Learning and Applications*. 117–122.

[5] Erkin Demirtas and Mykola Pechenizkiy. 2013. Cross-Lingual Polarity Detection with Machine Translation. In *Proceedings of the Second International Workshop on Issues of Sentiment Discovery and Opinion Mining* (Chicago, Illinois) *(WISDOM '13)*. Association for Computing Machinery, New York, NY, USA, Article 9, 8 pages. https://doi.org/10.1145/2502069.2502078

[6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. 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)*. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. https://doi.org/10.18653/v1/N19-1423

[7] A. Hayran and M. Sert. 2017. Analysis on Microblog Data based on Word Embedding and Fusion Techniques. In *2017 25th Signal Processing and Communication Application Conference (SIU)*. tbd. https://doi.org/tbd

[8] Mandar Joshi, Omer Levy, Luke Zettlemoyer, and Daniel Weld. 2019. BERT for Coreference Resolution: Baselines and Analysis. 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)*. Association for Computational Linguistics, Hong Kong, China, 5803–5808. https://doi.org/10.18653/v1/D19-1588

[9] Dilek Küçük, Guillaume Jacquet, and Ralf Steinberger. 2014. Named Entity Recognition on Turkish Tweets. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*. European Language Resources Association (ELRA), Reykjavik, Iceland, 450–454. http://www.lrec-conf.org/proceedings/lrec2014/pdf/380_Paper.pdf

[10] Deniz Kılınç, Akın Özçift, Fatma Bozyigit, Pelin Yıldırım, Fatih Yücalar, and Emin Borandag. 2017. TTC-3600: A new benchmark dataset for Turkish text categorization. *Journal of Information Science* 43, 2 (2017), 174–185. https://doi.org/10.1177/0165551515620551

[11] Guillaume Lample and Alexis Conneau. 2019. Cross-lingual Language Model Pretraining. *CoRR* abs/1901.07291 (2019). arXiv:1901.07291 http://arxiv.org/abs/1901.07291

[12] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. ALBERT: A Lite BERT for Self-supervised Learning of Language Representations. *ArXiv* abs/1909.11942 (2020).

[13] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *CoRR* abs/1907.11692 (2019). arXiv:1907.11692 http://arxiv.org/abs/1907.11692

[14] Tomas Mikolov, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013. Efficient Estimation of Word Representations in Vector Space. *CoRR* abs/1301.3781 (2013).

[15] Xiaoman Pan, Boliang Zhang, Jonathan May, Joel Nothman, Kevin Knight, and Heng Ji. 2017. Cross-lingual Name Tagging and Linking for 282 Languages. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Vancouver, Canada, 1946–1958. https://doi.org/10.18653/v1/P17-1178

[16] Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global Vectors for Word Representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Doha, Qatar, 1532–1543. https://doi.org/10.3115/v1/D14-1162

[17] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. *CoRR* abs/1802.05365 (2018). arXiv:1802.05365 http://arxiv.org/abs/1802.05365

[18] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language Models are Unsupervised Multitask Learners.

[19] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *ArXiv* abs/1910.10683 (2019).

[20] Stefan Schweter. 2020. BERTurk - BERT models for Turkish. https://doi.org/10.5281/zenodo.3770924

[21] Gökhan Akın Şeker and Gülşen Eryiğit. 2012. Initial explorations on using CRFs for Turkish named entity recognition. In *Proceedings of COLING 2012*. 2459–2474.

[22] Gokhan Tur, Dilek Hakkani-Tur, and Kemal Oflazer. 2003. A statistical information extraction system for Turkish. *Natural Language Engineering* 9 (06 2003), 181 – 210. https://doi.org/10.1017/S135132490200284X

[23] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. *CoRR* abs/1706.03762 (2017). arXiv:1706.03762 http://arxiv.org/abs/1706.03762

[24] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*. Association for Computational Linguistics, Brussels, Belgium, 353–355. https://doi.org/10.18653/v1/W18-5446

[25] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2016. Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. *CoRR* abs/1609.08144 (2016). arXiv:1609.08144 http://arxiv.org/abs/1609.08144

[26] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc Le. 2019. XLNet: Generalized Autoregressive Pretraining for Language Understanding. (06 2019).

[27] Reyyan Yeniterzi. 2011. Exploiting Morphology in Turkish Named Entity Recognition System. In *Proceedings of the ACL 2011 Student Session*. Association for Computational Linguistics, Portland, OR, USA, 105–110. https://www.aclweb.org/anthology/P11-3019

[28] Reyyan Yeniterzi, Gökhan Tür, and Kemal Oflazer. 2018. *Turkish Named-Entity Recognition.* Springer International Publishing, Cham, 115–132. https://doi.org/10.1007/978-3-319-90165-7_6

[29] Savaş Yildirim. 2014. A Knowledge-Poor Approach to Turkish Text Categorization. In *Proceedings of the 15th International Conference on Computational Linguistics and Intelligent Text Processing - Volume 8404* (Kathmandu, Nepal) *(CICLing 2014)*. Springer-Verlag, Berlin, Heidelberg, 428–440. https://doi.org/10.1007/978-3-642-54903-8_36

[30] Savaş Yildirim. 2020. *Comparing Deep Neural Networks to Traditional Models for Sentiment Analysis in Turkish Language*. Springer Singapore, Singapore, 311–319.

[31] Zhuosheng Zhang, Hai Zhao, and Rui Wang. 2020. Machine Reading Comprehension: The Role of Contextualized Language Models and Beyond. *ArXiv* abs/2005.06249 (2020).

[32] Gökhan Şeker and Gülşen Eryiğit. 2017. Extending a CRF-based named entity recognition model for Turkish well formed text and user generated content1. *Semantic Web* 8 (01 2017), 1–18. https://doi.org/10.3233/SW-170253