
ENHANCING DEEP NEURAL NETWORKS WITH MORPHOLOGICAL INFORMATION

A PREPRINT

Matej Klemen, Luka Krsnik, Marko Robnik-Šikonja

University of Ljubljana, Faculty of Computer and Information Science

Večna pot 113, Ljubljana, Slovenia

mk3141@student.uni-lj.si, krsnik.luka92@gmail.com, marko.robnik@fri.uni-lj.si

March 3, 2022

ABSTRACT

Deep learning approaches are superior in natural language processing due to their ability to extract informative features and patterns from languages. The two most successful neural architectures are LSTM and transformers, used in large pretrained language models such as BERT. While cross-lingual approaches are on the rise, most current natural language processing techniques are designed and applied to English, and less-resourced languages are lagging behind. In morphologically rich languages, information is conveyed through morphology, e.g., through affixes modifying stems of words. The existing neural approaches do not explicitly use the information on word morphology. We analyse the effect of adding morphological features to LSTM and BERT models. As a testbed, we use three tasks available in many less-resourced languages: named entity recognition (NER), dependency parsing (DP), and comment filtering (CF). We construct baselines involving LSTM and BERT models, which we adjust by adding additional input in the form of part of speech (POS) tags and universal features. We compare the models across several languages from different language families. Our results suggest that adding morphological features has mixed effects depending on the quality of features and the task. The features improve the performance of LSTM-based models on the NER and DP tasks, while they do not benefit the performance on the CF task. For BERT-based models, the added morphological features only improve the performance on DP when they are of high quality (i.e. manually checked) while not showing any practical improvement when they are predicted. Even for high-quality features, the improvements are less pronounced in language-specific BERT variants compared to massively multilingual BERT models. As in NER and CF datasets manually checked features are not available, we only experiment with predicted features and find that they do not cause any practical improvement in performance.

Keywords Deep learning · Natural language processing · Morphologically rich languages · Transformers · Morphological additions

1 Introduction

The use of deep learning for processing natural language is becoming a standard, with excellent results in a diverse range of tasks. Two state-of-the-art neural architectures for text-related modeling are long short-term memory (LSTM) networks [Hochreiter and Schmidhuber, 1997] and transformers [Vaswani et al., 2017]. LSTMs are recurrent neural networks that sequentially process text one token at a time, building up its internal representation in hidden states of the network. Due to the recurrent nature of LSTM, which degrades the efficiency of parallel processing, and improvements in performance, models based on the transformer architecture are gradually replacing LSTMs across many tasks. Transformers can process the text in parallel, using self-attention and positional embeddings to model the sequential nature of the text.

A common trend in using transformers is to pre-train them on large monolingual corpora with a general-purpose objective and then fine-tune them with a more specific objective, such as text classification. For example, the BERT (Bidirectional Encoder Representations from Transformers) architecture [Devlin et al., 2019] uses transformers and is pretrained with masked language modelling and order of sentences prediction tasks to build a general language understanding model. During the fine-tuning for a specific downstream task, additional layers are added to the BERT model, and the model is trained on task-specific data to capture the specific knowledge required to perform the task.

Most of the research in the natural language processing (NLP) area focuses on English, ignoring the fact that English is specific in terms of the low amount of information expressed through morphology (English is a so-called analytical language). In our work, we adapt modern deep neural networks, namely LSTM and BERT, for several morphologically rich languages by explicitly including the morphological information. The languages we analyse contain rich information about grammatical relations in the morphology of words instead of in particles or relative positions of words (as is the case in English). For comparison, we also evaluate our adaptations on English. Although previous research has shown that the state of the art methods such as BERT already capture some information contained in the morphology [Pires et al., 2019, Edmiston, 2020, Mikhailov et al., 2021], this investigation is commonly done by analysing the internals, for example with probing. Probing studies examine whether a property is encoded inside a model, but not necessarily used. In contrast, we present methods which combine BERT with separately encoded morphological properties: universal part of speech tags (UPOS tags) and universal features (grammatical gender, tense, conjugation, declination, etc.). We evaluate them on three downstream tasks: named entity recognition (NER), dependency parsing (DP), and comment filtering (CF), and observe whether the additional information benefits the models. If it does, the BERT models use the provided additional information, meaning that they do not fully capture it in pretraining. We perform similar experiments on LSTM networks and compare the results for both architectures.

Besides English, we analyse 10 more languages in NER, 15 in DP, and 5 in CF task. The choice of languages covers different language families but is also determined by the availability of resources and our limited computational resources. We describe the data in more detail in Section 3.

Our experiments show that the addition of morphological features has mixed effects depending on the task. Across the tasks where the added morphological features improve the performance, we show that 1) they benefit the LSTM-based models even if the features are noisy and 2) they benefit the BERT-based models only when the features are of high quality (i.e. human checked), suggesting that BERT models already capture the morphology of the language. We see room for improvement for large pretrained models either in designing pretraining objectives that can capture morphological properties or when high-quality features are available (rare in practice).

The remainder of this paper is structured as follows. In Section 2, we present different attempts to use morphological information in the three evaluation tasks and an overview of works studying the linguistic knowledge within neural networks. In Section 3, we describe the used datasets and their properties. In Section 4, we present the baseline models and models with additional morphological information, whose performance we discuss in Section 5. Finally, we summarize our work and present directions for further research in Section 6.

2 Related work

This section reviews the related work on the use of morphological information within the three evaluation tasks, mainly focusing on neural approaches. We split the review into four parts, one for each of the three evaluation tasks, followed by the works that study the linguistic knowledge contained within neural networks.

2.1 Morphological features in NER

Recent advances in NER are mostly based on deep neural networks. A common approach to NER is to represent the input text with word embeddings, followed by several neural layers to obtain the named entity label for each word on the output. In one of the earlier approaches, Collobert and Weston [2008] propose an architecture that is jointly trained on six different tasks, including NER, and show that this transfer learning approach generalizes better than networks trained on individual tasks due to learning a joint representation of tasks. While the authors use Time Delay Neural Network (convolutional) layers [Waibel et al., 1989] to model dependencies in the input, in subsequent works, various recurrent neural networks, such as LSTMs [Hochreiter and Schmidhuber, 1997], are commonly used. For example, Huang et al. [2015] show that using unidirectional and bidirectional LSTM layers for NER, results in comparable or better performance than approaches using convolutional layers.

In addition to word embeddings, these systems often use hand-crafted features, such as character-based features or affixes. The work of dos Santos and Guimarães [2015] outlines the inconvenience of constructing such features and proposes their automatic extraction using character embeddings and a convolutional layer. Authors combine the

character-level features with word-level features to obtain competitive or improved results on Spanish and Portuguese NER. Multiple authors [Kuru et al., 2016, Lample et al., 2016, Yang et al., 2016] confirm the effectiveness of character embeddings and show that recurrent layers can be used to process them instead of convolutional layers.

While sequence modelling on the character level can already encode morphological information, several authors show that the performance of neural networks on the NER task can be improved by including additional information about the morphological properties of the text. Straková et al. [2016] present a Czech NER system that surpasses the previous best system using only form, lemma and POS tag embeddings. Similarly, Güngör et al. [2017] show that using morphological embeddings in addition to character and word embeddings improves performance on Turkish and Czech NER, while Simeonova et al. [2019] show that including additional morphological and POS features improves performance on Bulgarian NER. Güngör et al. [2019] extend the study of Güngör et al. [2017] to three additional languages: Hungarian, Finnish and Spanish.

The influence of POS tags and morphological information on the NER performance of BERT models is less studied. Nguyen and Nguyen [2021] present a multi-task learning model which is jointly trained for POS tagging, NER, and DP. Their multi-task learning approach can be seen as an implicit injection of additional POS tags into BERT. The system trained on multiple tasks outperforms the respective single-task baselines, indicating that the additional information is beneficial. More explicitly, Mohseni and Tebbifakhr [2019] use morphological analysis as a preprocessing step to split the words into lemmas and affixes before passing them into a BERT model. Their system achieved first place in the NER shared task organized as part of the Workshop on NLP solutions for under-resourced languages 2019 [Taghizadeh et al., 2019]. However, they do not ablate the morphological analysis component, so it is unclear exactly how helpful it is in terms of the NER performance.

Our work extends the literature that studies the influence of POS tags and morphological information on the NER performance. For LSTM models, some works already explore this impact, though they are typically limited to one or a few languages. At the same time, we perform a study on a larger pool of languages from different language families. For BERT, we are not aware of any previous work that studies the influence of explicitly including POS tags or morphological information on the performance of the NER task. The existing work either adds the information implicitly [Nguyen and Nguyen, 2021] or does not show the extent to which the additional information is useful [Mohseni and Tebbifakhr, 2019]. In addition, existing analyses are limited to a single language. In contrast, our work explicitly adds POS tags or morphological information and shows how it affects the downstream (NER) performance on a larger pool of languages.

2.2 Morphological features in dependency parsing

Similarly to NER, recent progress in DP is dominated by neural approaches. Existing approaches introduce neural components into either a transition-based [Yamada and Matsumoto, 2003, Nivre, 2003] or a graph-based parser [McDonald et al., 2005]. Some works do not fall into either category, e.g., they treat DP as a sequence-to-sequence task [Li et al., 2018]. The two categories differ in how dependency trees are produced from the output of prediction models. In the transition-based approach, a model is trained to predict a sequence of parsing actions that produce a valid dependency tree. In contrast, in the graph-based approach, a model is used to score candidate dependency trees via the sum of scores of their substructures (e.g., arcs).

One of the earlier successful approaches to neural DP was presented by Chen and Manning [2014], who replaced the commonly used sparse features with dense embeddings of words, POS tags and arc labels, in combination with the transition-based parser. This approach improved both the accuracy and parse speed. Pei et al. [2015] introduced a similar approach to graph-based parsers. Later approaches improve upon the earlier methods by automatically extracting more information that guides the parsing, e.g., researchers use LSTM networks to inject context into local embeddings [Kiperwasser and Goldberg, 2016], apply contextual word embeddings [Kulmizev et al., 2019], or train graph neural networks [Ji et al., 2019].

The use of morphological features in DP, especially for morphologically rich languages, is common and predates neural approaches. For example, Marton et al. [2010] study the contribution of morphological features for Arabic, Seeker and Kuhn [2011] for German, Kapočiūtė-Dzikiene et al. [2013] for Lithuanian, Khallash et al. [2013] for Persian, etc. The majority of such works report improved results after adding morphological information. As our focus is on neural approaches, we mostly omit pre-neural approaches. Still, we note that the research area is extensive and has also been the topic of workshops such as the Workshop on statistical parsing of morphologically rich languages (Seddah et al. [2010]). The largely positive results have motivated authors to continue adding morphological information also to neural systems, which already automatically learn features and may pick up this information. For example, Chen and Manning [2014] note that POS tag embeddings contribute to the strong performance of their neural system on English and Chinese, and Özteş et al. [2018] note the usefulness of morphological embeddings for multiple agglutinative

languages. Dozat et al. [2017] report a similar trend in their work submitted to the CoNLL 2017 shared task [Hajič and Zeman, 2017]. They emphasize that POS tags are helpful, but only if produced using a sufficiently accurate POS tagger.

Similarly as in NER, character embeddings improve the accuracy of LSTM-based dependency parsers. We use the term “LSTM-based” to refer to models that include an LSTM neural network (as opposed to a more recent transformer neural network [Vaswani et al., 2017]). Several authors [Ballesteros et al., 2015, Dozat et al., 2017, Lhoneux et al., 2017] report that the use of character-based embeddings results in an improvement in the DP performance and can act as an approximate replacement for additional morphological information. Whether the embeddings present an approximate or complete replacement of morphological information is not entirely certain: Vania et al. [2018] show that models using character embeddings can still benefit from additional inclusion of morphological features. In contrast, Anderson and Gómez-Rodríguez [2020] report that the addition of POS tag embeddings does not further help a parser using character embeddings unless the POS tags are of a practically unrealistic quality.

Multi-task learning approaches, where a model is jointly trained for dependency parsing and another task (such as POS- or morphological tagging), are also a common way to inject additional information into dependency parsers [Straka, 2018, Lim et al., 2018, Nguyen and Verspoor, 2018]. Such approaches are popular in BERT-based parsers [Kondratyuk and Straka, 2019, Zhou et al., 2020a, Lim et al., 2020, Grünwald et al., 2021] and seem to be more common than the alternative approach of using additional inputs in a single task DP system. In our work, we use the alternative approach and explicitly include POS tags and morphological features in the form of their additional embeddings. We test the effect of the additional information on a sizable pool of languages from diverse language families. We perform several experiments to provide additional insight, testing the same effect with longer training, noisy information, and language-specific BERT models.

2.3 Morphological features in comment filtering

The literature for the CF task covers multiple related tasks, such as hate speech, offensive speech, political trolling, detecting commercialism, etc. Recent approaches involve variants of deep neural networks, though standard machine learning approaches are still popular, as shown in the survey of Fortuna and Nunes [2018]. These approaches typically use features such as character n-grams, word n-grams, and sentiment of the sequence. Two examples are the works of Malmasi and Zampieri [2017], who classify hate speech in English tweets, and Van Hee et al. [2015], who classify different levels of cyber-bullying in Dutch posts on ask.fm social site. Scheffler et al. [2018] combine word embeddings with the features mentioned above to classify German tweets. They observe that combining both n-gram features and word embeddings brings only a small improvement over only using one of them. The effectiveness of using features describing syntactic dependencies for toxic comments classification on English Wikipedia comments is shown by Shtovba et al. [2019].

Neural architectures used include convolutional neural networks [Georgakopoulos et al., 2018] and LSTM networks [Gao and Huang, 2017, Miok et al., 2019], typically improving the performance over standard machine learning approaches. The CF topic has also been the focus of shared tasks on identification and categorization of offensive language [Zampieri et al., 2019] and multilingual offensive language identification [Zampieri et al., 2020]. The reports of these tasks show the prevalence and general success of large pretrained contextual models such as BERT, though, surprisingly, the best performing model for the subtask B of SemEval-2019 Task 6 was rule-based [Han et al., 2019].

2.4 Linguistic knowledge combined with neural networks

Large pretrained models such as BERT show superior performance across many tasks. Due to a lack of theoretical understanding of this success, many authors study how and to what extent BERT models can capture various information, including different linguistic properties. An overview of recent studies in this area, sometimes referred to as BERTology, is compiled by Rogers et al. [2020]. Two common approaches to study BERT are i) add additional properties to BERT models and observe the difference in performance on downstream tasks, ii) a technique called probing [Conneau et al., 2018], where the BERT model is trained (fine-tuned) to predict a studied property. As we have noted examples of i) in previous sections, we focus on the probing attempts here.

For example, Jawahar et al. [2019] investigate what type of information is learned in different layers of the BERT English model and find that it captures surface features in lower layers, syntactic features in middle layers, and semantic features in higher layers. Similarly, Lin et al. [2019] find that BERT encodes positional information about tokens in lower layers and then builds increasingly abstract hierarchical features in higher layers. Tenney et al. [2019] use probing to quantify where different types of linguistic properties are stored inside BERT’s architecture and suggest that BERT implicitly learns the steps performed in classical (non-end-to-end) NLP pipeline. However, Elazar et al. [2021] point out possible flaws in the probing technique, suggesting amnesic probing as an alternative. They arrive at slightly

different conclusions about BERT layer importance; for example, they show that the POS information greatly affects the predictive performance in upper layers.

Probing studies for morphological properties were conducted by Edmiston [2020] and Mikhailov et al. [2021]. Concretely, they train a classifier to predict morphological features based on hidden layers of BERT. Based on the achieved high performance, Edmiston [2020] argues that monolingual BERT models capture significant amounts of morphological information and partition their embedding space into linearly-separable regions, correlated with morphological properties. Mikhailov et al. [2021] extend this work to multiple languages, performing probing studies on multilingual BERT models.

3 Data

In this section, we describe the datasets used in our experiments separately for each of the three tasks: NER, DP, and CF.

3.1 Named entity recognition

In the NER experiments, we use datasets in 11 languages from different language families: Arabic, Chinese, Croatian, English, Estonian, Finnish, Korean, Latvian, Russian, Slovene and Swedish. The number of sentences and tags present in the datasets is shown in Table 1. The label sets used in datasets for different languages vary, meaning that some contain more fine-grained labels than others. To make results across different languages consistent, we use IOB encoded labels present in all datasets: location (B/I-LOC), organization (B/I-ORG), person (B/I-PER), and “no entity” (O). We convert all other labels to the “no entity” label (O).

Table 1: The collected datasets for NER task and their properties: the number of sentences and tagged words. We display the results for the languages using their ISO 639-2 three letter code, provided in the “Code” column.

Language	Code	Dataset	Sentences	Tags
Arabic	ARA	ANERCorp [Benajiba et al., 2007]	5005	14876
Chinese	ZHO	MSRA [Levow, 2006]	48441	261940
Croatian	HRV	hr500k [Ljubešić et al., 2018]	24794	28902
English	ENG	CoNLL-2003 NER [Tjong Kim Sang and De Meulder, 2003]	20744	43979
Estonian	EST	Estonian NER corpus [Tkachenko et al., 2013]	14287	20965
Finnish	FIN	FiNER data [Ruokolainen et al., 2019]	14484	16833
Korean	KOR	KMOU NER	3659	6635
Latvian	LAV	LV Tagger train data [Paikens et al., 2012]	9903	11599
Russian	RUS	factRuEval-2016 [Starostin et al., 2016]	4907	9666
Slovene ¹	SLV	ssj500k [Krek et al., 2019]	9489	9440
Swedish	SWE	Swedish NER	9369	7292

3.2 Dependency parsing

To test morphological neural networks on the DP task, we use datasets in 16 languages from different language families: Arabic, Chinese, Croatian, English, Estonian, Finnish, Hebrew, Hungarian, Korean, Latvian, Lithuanian, Persian, Russian, Slovene, Swedish, and Turkish. We use the datasets from the Universal Dependencies [Nivre et al., 2020], which contain a collection of texts annotated with UPOS tags, XPOS tags (fine-grained POS), universal features, and syntactic dependencies. We provide the summary of the datasets in Table 2. The splits we use are predefined by the authors of the datasets. While most of the annotations are manually verified, the universal features in the Chinese and English datasets and the UPOS tags and universal features in the Turkish dataset are only partially manually verified. For Korean, the dataset does not contain universal feature annotations.

3.3 Comment filtering

While comparable datasets exist across different languages for the NER and DP task, no such standard datasets exist for the CF task. For that reason, in our experiments on CF, we select languages for which adequate datasets exist, i.e. large, of sufficient quality, and reasonably balanced across classes. We provide a summary of the used datasets in Table 3.

¹The Slovene ssj500k originally contains more sentences, but only 9489 are annotated with named entities.

Table 2: Dependency parsing datasets and their properties: the treebank, number of tokens, number of sentences, and the information about the size of splits. We display the dataset information using the language ISO 639-2 three-letter code provided in the ‘‘Code’’ column.

Language	Code	Treebank	Tokens	Sentences	Train	Validation	Test
Arabic	ARA	PADT	282384	7664	6075	909	680
Chinese	ZHO	GSD	123291	4997	3997	500	500
Croatian	HRV	SET	199409	9010	6914	960	1136
English	ENG	EWT	254855	16622	12543	2002	2077
Estonian	EST	EDT	438171	30972	24633	3125	3214
Finnish	FIN	TDT	202697	15135	12216	1364	1555
Hebrew	HEB	HTB	161411	6216	5241	484	491
Hungarian	HUN	Szeged	42032	1800	910	441	449
Korean	KOR	Kaist	350090	27363	23010	2066	2287
Latvian	LAV	LVTB	220536	13643	10156	1664	1823
Lithuanian	LIT	ALKSNIS	70051	3642	2341	617	684
Persian	FAS	PerDT	501776	29107	26196	1456	1455
Russian	RUS	GSD	98000	5030	3850	579	601
Slovene	SLV	SSJ	140670	8000	6478	734	788
Swedish	SWE	Talbanken	96858	6026	4303	504	1219
Turkish	TUR	BOUN	122383	9761	7803	979	979

For English experiments, we use a subset of toxic comments from Wikipedia’s talk page edits used in Jigsaw’s toxic comment classification challenge [Wulczyn et al., 2017]. The comments are annotated with six possible labels: toxic, severe toxic, obscene language, threats, insults, and identity hate (making a total of six binary target variables). We extracted comments from four categories: toxic, severe toxic, threats, and identity hate, a total of 21, 541 instances. We randomly chose the same amount of comments that do not fall in any of the mentioned categories, obtaining the final dataset of 43, 082 instances, using 60% randomly selected examples as the training, 20% as the validation, and 20% as the test set.

For Korean experiments, we use a dataset of comments from a Korean news platform [Moon et al., 2020], annotated as offensive, hateful or clean. We group the offensive and hateful examples to produce a binary classification task. As the test set labels are private, we instead use the predefined validation set as the test set and set aside 20% of the training set as the new validation set.

For Slovene experiments, we use the IMSyPP-sl dataset [Evkoski et al., 2021], containing tweets annotated for fine-grained hate speech. The tweets are annotated with four possible labels: appropriate (i.e. not offensive), inappropriate, offensive, or violent. Each tweet is annotated twice, and we only keep a subset for which both labels agree. To produce a binary classification task, we group the tweets labelled as inappropriate, offensive, or violent into a single category. We use the predefined split into a training and test set, and additionally, remove 20% of examples from the training set for use in the validation set.

For Arabic, Greek, and Turkish experiments, we use datasets from the OffensEval 2020 shared task on multilingual offensive language identification [Zampieri et al., 2020]. The datasets are composed of tweets annotated for offensive language: a tweet is either deemed offensive or not offensive. We use the predefined splits into training and test sets provided by the authors. We randomly remove 20% of the examples from the original training sets to create the validation sets. As the Turkish training set proved to be too heavily imbalanced, we decided to randomly remove half of the unoffensive examples from it before creating the validation set.

4 Neural networks with morphological features

This section describes the architectures of neural networks used in our experiments. Their common property is that we enhance standard word embeddings based inputs with embeddings of morphological features. We work with recent successful neural network architectures, LSTMs and transformers, i.e. BERT models. A detailed description of architectures is available in the following subsections, separately for each evaluation task. We describe the baseline architecture and the enhanced one for each task and architecture.

Table 3: Comment filtering datasets and their properties: number of examples, size of the split and class distribution inside subsets.

Language	Dataset	Examples	Train	Validation	Test
Arabic	OffensEval 2020 [Zampieri et al., 2020]	9959	6370	1597	1992
English	Jigsaw toxic comments [Wulczyn et al., 2017]	43082	25848	8616	8618
Greek	OffensEval 2020 [Zampieri et al., 2020]	10287	6994	1749	1544
Korean	Korean hate speech dataset [Moon et al., 2020]	8367	6316	1580	471
Slovene	IMSyPP-sl[Evkoski et al., 2021]	47538	31676	7919	7943
Turkish	OffensEval 2020 [Zampieri et al., 2020]	22470	15154	3789	3527

4.1 Named entity recognition models

In the NER task, we use two baseline neural networks (LSTM and BERT) and the same two models with additional morphological information: POS tag embeddings and universal feature embeddings. The baseline models and their enhancements are displayed in Figure 1.

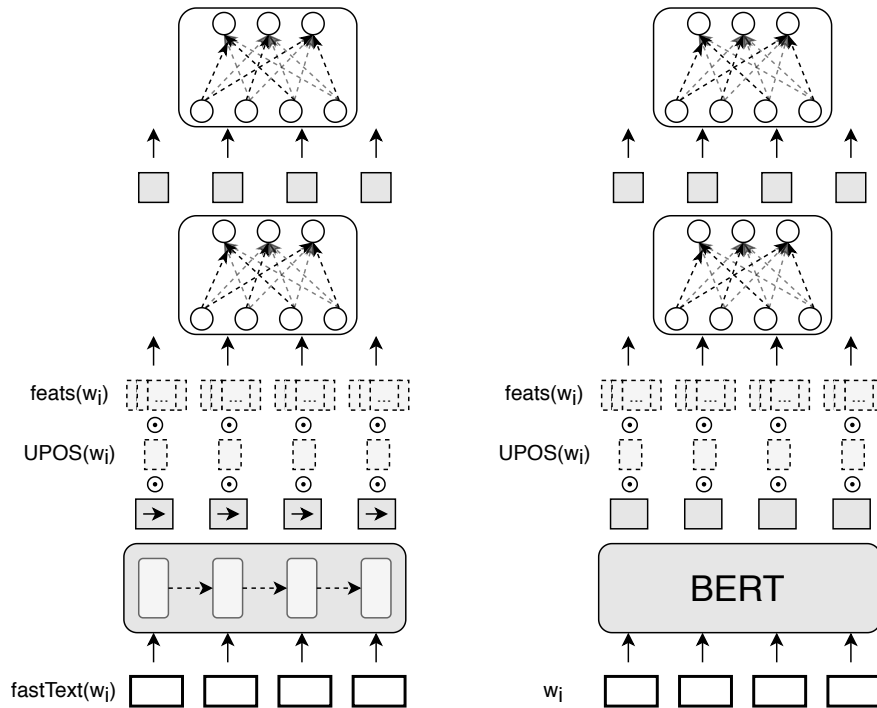


Figure 1: The baseline LSTM-based (left) and BERT-based (right) models for the NER task, along with our modifications with morphological information. The dotted border of POS vectors and morphological feature vectors (feats) marks that their use is optional and varies across experiments. The \odot symbol between layers represents the concatenation operation. The w_i symbol stands for token i ; in case of LSTM, tokens enter the model sequentially, and we show the unrolled network, while BERT processes all tokens simultaneously.

The first baseline model (left-hand side of Figure 1) is a unidirectional (left-to-right) LSTM model, which takes as an input a sequence of tokens, embedded using 300-dimensional fastText embeddings [Bojanowski et al., 2017]. These embeddings are particularly suitable for morphologically rich languages as they work with subword inputs². For each input token, its LSTM hidden state is extracted and passed through the linear layer to compute its tag probabilities.

The second baseline model (right-hand side of Figure 1) is the cased multilingual BERT base model (bert-base-multilingual-cased). In our experiments, we follow the sequence tagging approach suggested by the authors of

²The precomputed embeddings are available at <https://fasttext.cc/docs/en/crawl-vectors.html>.

BERT [Devlin et al., 2019]. The input sequence is prepended with a special token [CLS] and passed through the BERT model. The output of the last BERT hidden layer is passed through the linear layer to obtain the predictions for NER tags.

Both baseline models (LSTM and BERT) are enhanced with the same morphological information: POS tag embeddings and universal feature embeddings for each input token. We embedded the POS tags using 5-dimensional embeddings. For each of the 23 universal features used (we omitted the *Typo* feature, as the version of the used tagger did not annotate this feature), we independently constructed 3-dimensional embeddings, meaning that we obtained a 69-dimensional universal embedding. We embedded the features independently due to a large number of their combinations and treated them equally in the DP and CF experiments. We selected the size of the embeddings based on the results of preliminary experiments on Slovene and Estonian languages. We automatically obtained the POS tags and morphological features using the Stanza system [Qi et al., 2020]. In the enhanced architectures, we included another linear layer before the final linear classification layer to model possible interactions.

4.2 Dependency parsing models

As the baseline model in the DP task, we use the deep biaffine graph-based dependency parser [Dozat and Manning, 2016]. The enhancements with the morphological information are at the input level. The baseline model and its enhancements are shown in Figure 2.

The baseline parser combines a multi-layer bidirectional LSTM network with a biaffine attention mechanism to jointly optimize prediction of arcs and arc labels. We leave the majority of baseline architectural hyperparameters at values described in the original paper (3-layer bidirectional LSTM with 100-dimensional input word embeddings and the hidden state size of 400).

In our experiments, we concatenate the non-contextual word embeddings with various types of additional information. The first additional input is contextual word embeddings, which we obtain either by using the hidden states of an additional single-layer unidirectional LSTM or by using a learned linear combination of all hidden states of an uncased multilingual BERT base model (bert-base-multilingual-uncased). To check whether the results depend on using cased or uncased BERT model, we rerun experiments for a small sample of languages with the cased version, using identical hyperparameters, and present the results in Appendix B. The conclusions drawn from both types are similar; the improvements of enhanced models are statistically insignificant for only one out of the eight sampled languages when using a cased BERT model.

Although the LSTM layers are already present in the baseline parser, we include an additional LSTM layer at the input level to explicitly encode the context, keeping the experimental settings similar across our three evaluation tasks, i.e. we have one setting with added LSTM and one setting with added BERT. The second additional input is universal POS embeddings (UPOS), and the third is universal feature embeddings (feats). These embeddings are concatenated separately for each token of the sentences. The size of the additional LSTM layer, POS tag embeddings and universal feature embedding are treated as tunable hyperparameters. As the baseline input embeddings, we use pre-trained 100-dimensional fastText embeddings, which we obtain by reducing the dimensionality of publicly available 300-dimensional vectors with fastText’s built-in dimensionality reduction tool.

In DP experiments, we use POS tags and morphological features of two origins. The first source is human annotations provided in the used datasets. The second source of morphological information is predictions of Stanza models [Qi et al., 2020]. These two origins are used to assess the quality of morphological information; namely, we check if manual human annotations provide any benefit compared to automatically determined POS tags and features.

4.3 Comment filtering models

We add additional morphological information to standard LSTM and BERT models in the CF evaluation task. The baseline and enhanced models for this task are similar to those in the NER evaluation, though we operate at the sequence level here instead of the token level in the NER task. The architecture of models is shown in Figure 3. As baselines, we take a single layer unidirectional LSTM network (the top part of Figure 3) and the multilingual base uncased BERT model (bert-base-multilingual-uncased; the bottom part of Figure 3). The difference in the used BERT dialect (*uncased* as opposed to the *cased* in the NER task) is due to better performance detected in preliminary experiments.

In the LSTM baseline model, the words of the input sequence are embedded using pre-trained 300-dimensional fastText embeddings. As the representation of the whole sequence, we take the output of the last hidden state, which then passes through the linear layer to obtain the prediction scores. In the BERT baseline model, we take the sequence classification approach suggested by the authors of BERT. The input sequence is prepended with the special [CLS] token and passed

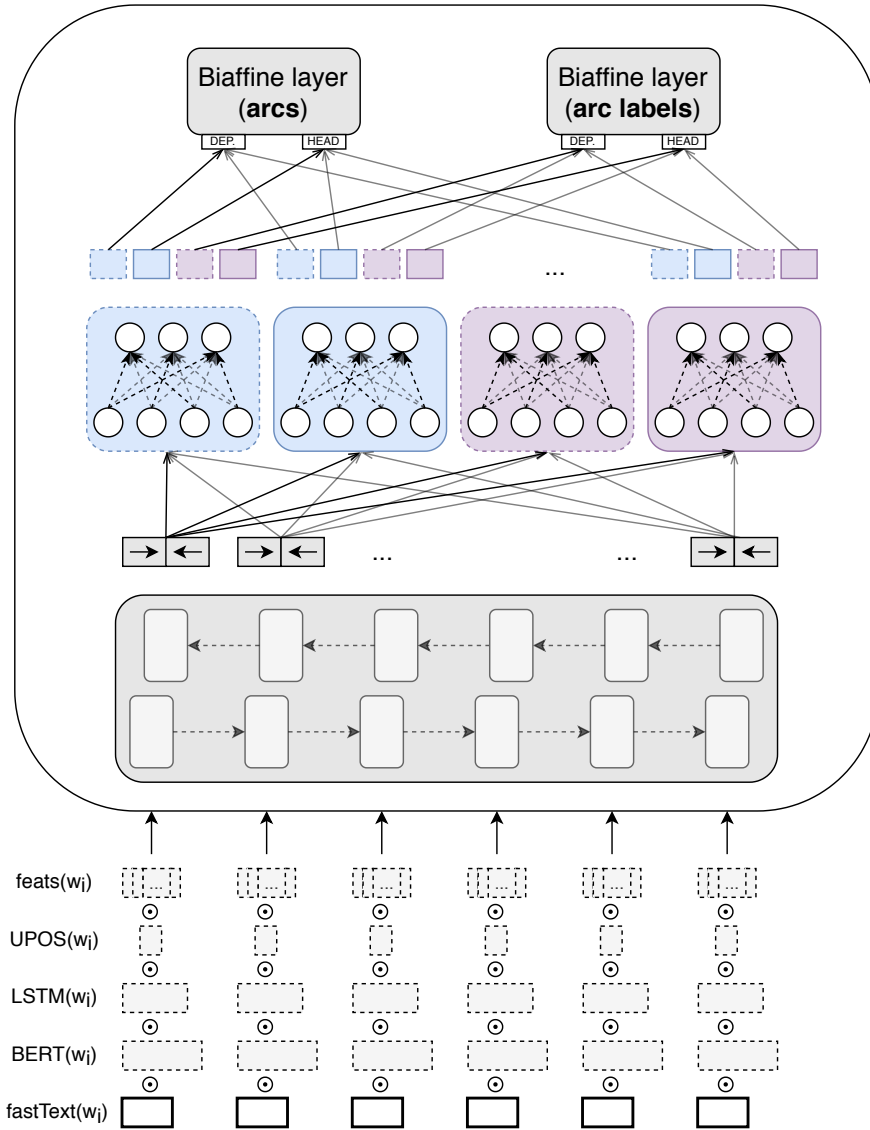


Figure 2: The deep biaffine graph-based dependency parser along with our enhancements at the input level. The dotted border of input embedding vectors, POS vectors, and morphological features (feats) is optional and varies across experiments. The \odot symbol between layers represents the concatenation operation. The w_i symbol stands for token i ; tokens enter the LSTM model sequentially, and we show the unrolled network.

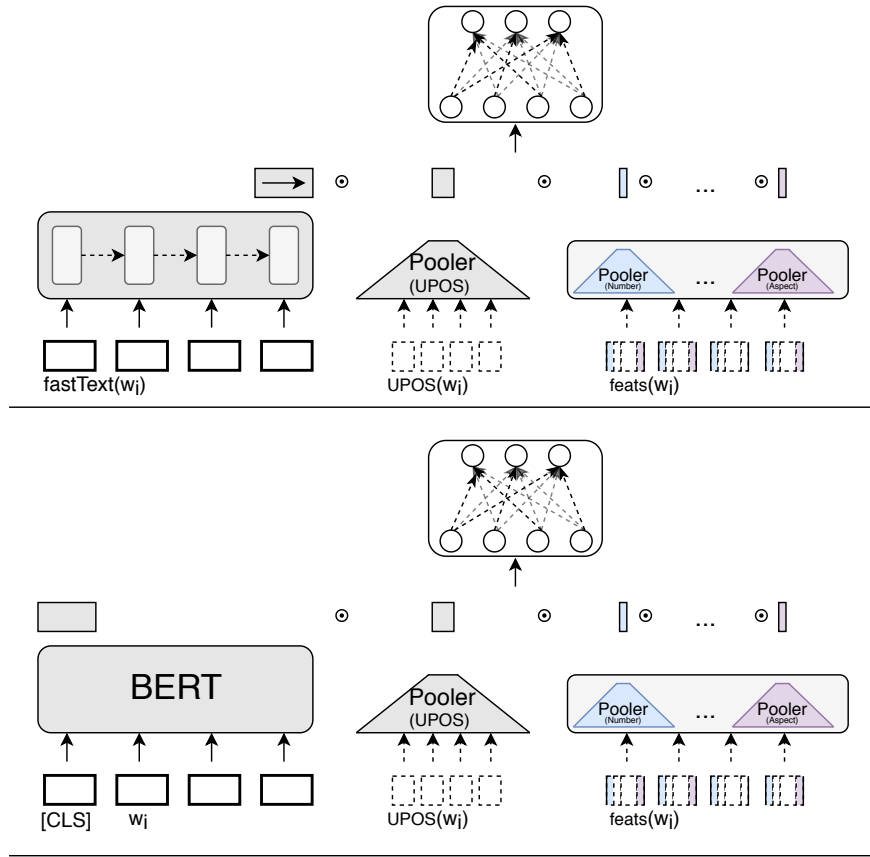


Figure 3: The baseline LSTM (top) and BERT (bottom) models for the CF task, along with our modifications with morphological information. The dotted border of UPOS vectors and morphological features (feats) marks that their use is optional and varies across experiments. The \odot symbol between layers represents the concatenation operation. The w_i symbol stands for token i ; in the case of LSTM, tokens enter the model sequentially, and we show the unrolled network, while BERT processes all tokens simultaneously.

through BERT. The sequence representation corresponds to the output of the last BERT hidden layer for the [CLS] token, which is passed through a linear layer to obtain the prediction scores.

We augment the baseline models with POS tags and universal feature embeddings, as in other tasks. We obtain the tags for each token separately using the Stanza system [Qi et al., 2020]. We combine the obtained embeddings using three different pooling mechanisms: mean, weighted combination, or LSTM pooling. Given the POS tag or universal feature embeddings, the mean pooling outputs the mean of all token embeddings. The weighted pooling outputs the weighted combination of token embeddings, and the LSTM pooling outputs the last hidden state obtained by passing the sequence of embeddings through the LSTM network. Both the embedding sizes and the type of pooling are treated as tunable hyperparameters. The coefficients of the weighted combination are learned by projecting the sequence embeddings into a sequence of independent dimension values, which are normalized with the softmax. This compresses the embedding sequence, establishes its fixed length, and allows different morphological properties to have a different impact on the sequence representation. This approach tests the contextual encoding of morphological properties. For example, we might learn that the adjectives are assigned a higher weight than other POS tags due to their higher emotional contents that often indicates insults.

5 Evaluation

In this section, we first present the evaluation scenario for the three evaluation tasks, followed by the results presented separately for each of the tasks. We end the section with additional experiments performed on the DP task. In the

additional experiments, we investigate the effect of additional morphological features in three situations where we tweak one aspect of the training procedure at a time: 1) we increase the maximum training time of our models by additional 5 epochs, 2) we replace the human-annotated features with the ones automatically predicted by machine learning models, and 3) we replace the embeddings from general multilingual BERT with those from more specific multilingual and monolingual BERT models.

5.1 Experimental settings

The experimental settings differ between the three evaluation tasks, so we describe them separately for each task, starting with the NER task and then the DP and CF tasks. One aspect that is common across all three tasks is that we skip the experiments involving universal features on Korean, as the features are not available for any of the bigger Korean Universal Dependencies corpora; therefore, the Stanza models cannot be used to predict them either.

5.1.1 Experimental settings for NER

For NER, we train each BERT model for 10 epochs and each LSTM model for 50 epochs. These parameters were determined during preliminary testing on the Slovene dataset. The selected numbers of epochs are chosen to balance the performance and training times of the models. All NER models are evaluated using 10-fold cross-validation.

We assess the performance of the models with the F_1 score, which is a harmonic mean of precision and recall measures. This measure is typically used in NER in a way that precision and recall are calculated separately for each of the three entity classes (location, organization, and person, but not “no entity”). We compute the weighted average over the class scores for each metric, using the frequencies of class values in the datasets as the weights. Ignoring the “no entity” (O) label is a standard approach in the NER evaluation and disregards words that are not annotated with any of the named entity tags, i.e. the assessment focuses on the named entities.

5.1.2 Experimental settings for dependency parsing

We train each DP model for a maximum of 10 epochs, using an early stopping tolerance of 5 epochs. All models are evaluated using predefined splits into training, validation and testing sets, determined by the respective treebank authors and maintainers (see Table 2). The final models, evaluated on the test set, are selected based on the maximum mean of unlabeled and labelled attachment scores (UAS and LAS) on the validation set. The UAS and LAS are standard accuracy metrics in DP. The UAS score is defined as the proportion of tokens assigned the correct syntactic head. In contrast, the LAS score is the proportion of tokens that are assigned the correct syntactic head as well as the dependency label [Jurafsky and Martin, 2009]. As both scores are strongly correlated in most models, we only report the LAS scores to make the presentation clearer.

5.1.3 Experimental settings for comment filtering

In our CF experiments, we train BERT models for a maximum of 10 epochs and LSTM models for a maximum of 50 epochs, using early stopping with the tolerance of 5 epochs.

We evaluate the models using the macro F_1 score, computed as the unweighted mean of the F_1 scores for each label. We use fixed training, validation and test sets, described in Table 3. As we have observed noticeable variance in the preliminary experiments, we report the mean metrics over five runs for each setting, along with the standard deviation.

5.2 Experimental results

In this section, we compare the results of baseline models with their enhancements using additional morphological information. We split the presentation into three parts, according to the evaluation task: NER, DP, and CF.

5.2.1 Results for the NER task

For the NER evaluation task, we present results of the baseline NER models and models enhanced with the POS tags and universal features, as introduced in Section 4.1. Table 4 shows the results for LSTM and BERT models for 11 languages. We compute the statistical significance of the differences between the baseline LSTM and BERT models and their best-performing counterparts with morphological additions. We use the Wilcoxon signed-rank test [Wilcoxon et al., 1970] and underline the statistically significant differences at $p = 0.01$ level.

The baseline models involving BERT outperform their LSTM counterparts across all languages by a large margin. When adding POS tags or universal features to the LSTM-based models, we observe an increase in performance over

Table 4: F_1 scores of different models on the NER task in 11 languages. The left part of the table shows the results for the LSTM models and the right part for the BERT models. The best scores for each language and neural architecture are marked with the bold typeface. Best scores for which the difference to the respective baseline is statistically significant are underlined.

(a) LSTM.					(b) BERT.			
lang.	+LSTM	+LSTM +UPOS	+LSTM +UPOS +feats	+LSTM +feats	+mBERT	+mBERT +UPOS	+mBERT +UPOS +feats	+mBERT +feats
ARA	0.690	0.690	0.689	0.691	0.821	0.816	0.817	0.821
ENG	0.890	0.890	0.884	0.886	0.948	0.946	0.948	0.947
EST	0.769	0.789	0.786	0.777	0.875	0.879	0.877	0.875
FIN	0.814	0.828	0.829	0.826	0.926	0.924	0.928	0.929
HRV	0.700	0.713	0.715	0.704	0.874	0.871	0.872	0.874
KOR	0.516	0.519	/	/	0.881	0.882	/	/
LAV	0.573	0.588	0.581	0.581	0.766	0.776	0.780	0.777
RUS	0.752	0.753	0.765	0.760	0.868	0.868	0.870	0.869
SLV	0.651	0.696	0.695	0.669	0.848	0.849	0.851	0.854
SWE	0.785	0.783	0.787	0.793	0.885	0.886	0.885	0.880
ZHO	0.787	0.786	0.784	0.781	0.930	0.930	0.930	0.930

the baselines for nine languages. For five (Arabic, Korean, Latvian, Russian, and Swedish), the increase in the F_1 score is not statistically significant; for three of them, the difference is under 0.01. For the remaining four languages (Croatian, Estonian, Finnish, and Slovene), the increase is statistically significant and ranges from 0.015% (Finnish and Croatian) to 0.045% (Slovene). The two languages for which we observe no improvement after adding morphological information to the LSTM-based models are English and Chinese.

In BERT-based models, the additional information does not make a practical difference. For all languages except Slovene and Latvian, the increase in F_1 values over the baseline is below 0.005. However, we note that the better results on Slovene are possibly a result of the optimistically high quality of the additional morphological information. The Stanza model with which we obtained the predictions for Slovene was trained for a different task, but on the same dataset as we use for NER. Still, the improvements are not statistically significant, so we do not explore the effect further.

5.2.2 Results for the dependency parsing task

For the DP evaluation task, we present LAS scores in Table 5. On the left, we show the results obtained with the parser containing additional LSTM embeddings, and on the right, with the parser containing additional BERT embeddings. For each, we report the results without and with additional POS tag and universal feature inputs (as introduced in Section 4.2). We also statistically test the differences in the scores between the best performing enhanced variants and their baselines without morphological additions. As the splits are fixed in the DP tasks, we use the Z-test for the equality of two proportions [Kanji, 2006] at $p = 0.01$ level. The null hypothesis is that the scores of compared models are equal. We underline the best result for languages where the null hypothesis can be rejected.

Similarly to NER, the models involving BERT embeddings outperform the baselines involving LSTM embeddings on all 16 languages by a large margin. The models with the added POS tags or universal features noticeably improve over baselines with only LSTM embeddings for all languages. The increase ranges between 2.44% (Russian) and 10.42% (Lithuanian). All compared differences between the LSTM baselines and the best enhanced variants are statistically significant at $p = 0.01$.

Contrary to the results observed on the NER task, adding morphological features to the models with BERT embeddings improves the performance scores for all languages. The increase ranges from 0.57% (Russian) and 4.45% (Chinese). The increase is not statistically significant only for one language (Russian). For the remaining 15 languages, the increase is over one per cent and statistically significant.

Interestingly, for some languages, the enhanced parsers with LSTM come close in terms of LAS to the baseline parsers with BERT. For two languages (Lithuanian and Latvian), they even achieve better LAS scores.

Table 5: LAS scores for different models on the DP task in different languages using (a) LSTM or (b) BERT contextual embeddings. The best scores for each language and each type of contextual embeddings are marked with the bold typeface. Statistically significant differences at $p = 0.01$ level in the best scores compared to the baseline of the same architecture are underlined.

lang.	(a) LSTM.				(b) BERT.			
	+LSTM	+LSTM +UPOS	+LSTM +UPOS +feats	+LSTM +feats	+mBERT	+mBERT +UPOS	+mBERT +UPOS +feats	+mBERT +feats
ARA	76.74	80.11	80.53	79.42	80.97	82.54	83.01	82.45
ENG	80.69	84.68	84.48	82.96	88.09	89.35	89.13	88.81
EST	76.24	81.36	83.46	81.00	84.30	86.72	87.19	85.72
FAS	84.42	87.59	87.32	85.95	89.06	90.80	90.65	89.90
FIN	78.60	83.06	83.97	81.16	86.07	87.32	87.99	87.07
HEB	82.29	85.49	84.92	83.84	87.18	88.38	88.77	87.63
HRV	79.59	81.95	82.84	82.00	86.37	87.63	87.44	86.81
HUN	69.58	75.19	75.93	74.29	77.07	80.04	79.96	77.05
KOR	74.46	82.35	/	/	85.46	87.59	/	/
LAV	78.08	82.75	84.27	81.81	83.02	85.43	86.16	85.19
LIT	63.65	70.45	74.07	72.40	72.62	74.99	77.04	76.33
RUS	79.24	80.88	81.68	81.64	85.99	86.24	86.56	86.37
SLV	83.03	88.75	90.32	88.38	91.67	92.90	93.40	92.74
SWE	78.97	82.65	83.17	80.31	86.67	87.80	88.38	87.25
TUR	66.85	68.62	69.46	67.51	70.60	72.99	72.58	71.58
ZHO	72.04	78.65	78.89	71.09	79.49	83.94	83.21	80.45

5.2.3 Results for the comment filtering task

Table 6 shows the CF results for LSTM and BERT baselines and their enhancements, described in Section 4.3. We compute the statistical significance of the differences between the baseline LSTM and BERT models and their best-performing counterparts with morphological additions. To do so, we use the unpaired t-test [Student, 1908] and check for statistical significance at $p = 0.01$ level, as in the other two tasks.

All BERT-based models outperform the LSTM-based models by a large margin in all languages. Adding POS tags and universal features to neural architectures of either type does not seem to significantly benefit their F_1 score in general. Although the mean scores of the best performing enhanced models are often higher, the difference is under 0.010 in most cases. The exceptions are Greek (+0.017) and Slovene (+0.010) for LSTM models, and Korean (+0.010) for BERT models, although the improvements are not statistically significant due to the noticeable standard deviation of the scores. On the other hand, the enhanced models perform practically equivalently or slightly worse on Arabic. As the performance scores do not statistically differ with the best set of hyperparameters, we do not further analyse the effect of different pooling types on the performance. However, we did not observe any pooling approach to perform best consistently.

5.3 Additional experiments

To further analyse the impact of different aspects of the proposed morphological enhancements, we conducted several studies on the DP task, where the datasets and evaluation settings allow many experiments. Similarly as in Section 5.2.2, we statistically evaluate differences in performance between the baseline model and the best enhancement using the Z-test for the equality of two proportions. In cases where the null hypothesis can be rejected at $p = 0.01$ level, we underline the respective compared score. We test the impact of additional training time (Section 5.3.1), quality of morphological information (Section 5.3.2), and different variants of BERT models (Section 5.3.3).

5.3.1 Additional training time

To test if the observed differences in performance are due to random variation in the training of models, which could be reduced with longer training times, we increase the maximum training time from 10 to 15 epochs. We show the results in Table 7.

Table 6: F_1 scores for baseline and enhanced models on the CF task in different languages. The left part of the table shows results for LSTM models, and the right part shows the results for BERT models. We report the mean and standard deviation over five runs. The highest mean score for each language is marked with the bold typeface separately for LSTM and BERT models. Note, however, that none of the improvements over the baselines is statistically significant at $p = 0.01$ level.

lang.	(a) LSTM.				(b) BERT.			
	+LSTM	+LSTM +UPOS	+LSTM +UPOS +feats	+LSTM +feats	+mBERT	+mBERT +UPOS	+mBERT +UPOS +feats	+mBERT +feats
ARA	0.752 (0.005)	0.748 (0.006)	0.749 (0.006)	0.749 (0.006)	0.843 (0.005)	0.843 (0.007)	0.840 (0.008)	0.842 (0.008)
ELL	0.646 (0.014)	0.647 (0.012)	0.663 (0.009)	0.652 (0.014)	0.820 (0.022)	0.821 (0.014)	0.823 (0.015)	0.817 (0.011)
ENG	0.881 (0.012)	0.882 (0.015)	0.869 (0.012)	0.883 (0.027)	0.922 (0.012)	0.923 (0.011)	0.923 (0.020)	0.922 (0.013)
KOR	0.661 (0.012)	0.661 (0.008)	/	/	0.722 (0.007)	0.732 (0.019)	/	/
SLV	0.640 (0.013)	0.650 (0.015)	0.650 (0.006)	0.649 (0.007)	0.782 (0.006)	0.784 (0.005)	0.785 (0.002)	0.782 (0.005)
TUR	0.671 (0.010)	0.676 (0.010)	0.669 (0.008)	0.676 (0.009)	0.756 (0.005)	0.762 (0.008)	0.765 (0.006)	0.762 (0.010)

We can observe that longer training times slightly increases the scores for all model variants, though their relative order stays the same. The models with added morphological features still achieve better results, so the performance increases do not seem to be the effect of random fluctuations in training due to the number of training steps. All improvements over baselines of the parsers using LSTM embeddings remain statistically significant. In contrast, for the parsers using BERT embeddings, the improvements for two languages are now no longer statistically significant (Croatian and Russian).

5.3.2 Quality of morphological information

In the second additional experiment, we evaluate the impact of the quality of morphological information. We replace the high-quality (human-annotated) POS tags and morphological features used in Section 4.2 with those predicted by machine learning models. In this way, we test a realistic setting where the morphological information is at least to a certain degree noisy. We obtain POS tags and morphological features from Stanza models prepared for the tested languages [Qi et al., 2020]. To avoid overly optimistic results, we use models that are not trained on the same datasets used in our DP experiments. This is possible for a subset of nine languages. We note the used Stanza models in Appendix A, together with the proportion of tokens, for which POS tags and *all* universal features are correctly predicted (i.e. their accuracy).

We show the results of DP models, trained with predicted morphological features in Table 8.

The general trend is that using predicted features results in much smaller (best case) performance increases, though some languages still see significant increases. For LSTM models, the increases range from 0.62% (Persian) to 2.41% (Finnish), with six out of nine being statistically significant. For BERT models, the increases are all statistically insignificant and range from -0.10% (i.e. decrease, Slovene) to 0.66% (Finnish). These results are consistent with the results of our NER experiments in Section 5.2.1, where we have no access to human-annotated features and find that noisy features only help LSTM-based models.

These results indicate that adding predicted morphological features to models with BERT embeddings might not be practically useful, since their quality needs to be very high. However, since human-annotated morphological features improve the performance on the DP task, this suggests that there could be room for improvement in BERT pre-training. It seems that the pre-training tasks of BERT (masked language modelling and next sentence prediction) do not fully capture the morphological information present in the language. However, it is unlikely that the models could capture all

Table 7: LAS scores achieved by models that are trained for up to 5 additional epochs (a maximum training time of 15 epochs instead of 10). The results for 10 epochs are presented in Table 5. Statistically significant differences in best scores at $p = 0.01$ level compared to the baselines of the same architecture are underlined.

(a) LSTM.					(b) BERT.			
lang.	+LSTM	+LSTM +UPOS	+LSTM +UPOS +feats	+LSTM +feats	+mBERT	+mBERT +UPOS	+mBERT +UPOS +feats	+mBERT +feats
ARA	78.00	80.56	81.13	79.68	81.75	82.77	83.43	82.95
ENG	81.47	84.74	85.26	83.58	87.77	89.56	89.49	88.89
EST	76.72	81.52	83.65	81.24	84.60	86.60	87.26	85.78
FAS	84.82	87.61	87.32	85.59	88.90	90.88	90.75	89.94
FIN	80.28	83.77	84.22	81.95	87.06	87.98	88.42	87.85
HEB	83.06	86.35	85.95	84.83	88.02	88.88	89.34	88.83
HRV	79.49	82.72	83.26	82.74	87.03	87.55	87.71	86.97
HUN	70.41	76.18	77.43	75.59	79.31	81.29	81.88	79.71
KOR	74.41	82.28	/	/	85.58	87.84	/	/
LAV	79.58	83.32	84.71	82.50	84.04	85.66	86.73	85.19
LIT	66.48	71.03	74.22	74.28	74.17	76.11	77.88	77.04
RUS	79.79	81.35	82.18	82.21	86.15	87.04	87.23	86.64
SLV	84.72	89.19	90.64	88.60	92.11	92.80	93.78	93.08
SWE	79.97	84.58	85.23	82.10	87.43	88.86	88.77	87.72
TUR	67.68	69.86	70.42	69.12	71.78	71.82	73.38	72.24
ZHO	72.94	80.45	79.60	72.69	81.29	84.19	84.01	80.85

Table 8: LAS scores achieved by DP models that are trained with predicted (noisy) instead of human-annotated morphological features. We provide the accuracy of predicted UPOS tags and universal features in Appendix A. Statistically significant differences in the best scores compared to baselines at $p = 0.01$ level are underlined.

(a) LSTM.					(b) BERT.			
lang.	+LSTM	+LSTM +UPOS	+LSTM +UPOS +feats	+LSTM +feats	+mBERT	+mBERT +UPOS	+mBERT +UPOS +feats	+mBERT +feats
ENG	80.69	81.36	81.02	81.22	88.09	88.33	88.05	87.74
EST	76.24	77.15	77.39	76.19	84.30	84.69	84.39	84.42
FIN	78.60	80.61	81.01	80.43	86.07	86.33	86.73	85.91
KOR	74.46	76.23	/	/	85.46	85.90	/	/
FAS	84.42	85.04	84.07	84.59	89.06	89.28	88.82	88.84
RUS	79.24	79.53	80.89	79.56	85.99	86.17	85.74	85.27
SLV	83.03	84.06	85.00	84.12	91.67	91.42	91.57	91.23
SWE	78.97	80.06	80.08	79.24	86.67	86.59	86.96	86.33
TUR	66.85	67.42	68.03	67.76	70.60	71.26	71.11	70.77

information present in the ground truth annotations, as humans can disambiguate the grammatical role of a word even where syncretism occurs.

5.3.3 Variants of BERT model

In the third additional experiment, we revert to using the ground truth morphological annotations but replace the embeddings obtained from the multilingual uncased BERT model with those obtained from more specific multilingual BERT models and monolingual BERT models. In experiments involving multilingual BERT models, we test the Croatian/Slovene/English CroSloEngual BERT, Finnish/Estonian/English FinEst BERT [Ulčar and Robnik-Šikonja, 2020], and Bulgarian/Czech/Polish/Russian Slavic BERT [Arhipov et al., 2019]. In experiments with monolingual BERT models, we use the Arabic bert-base-arabic [Safaya et al., 2020], English bert-base-cased [Devlin et al., 2019], Estonian ESTBert [Tanvir et al., 2021], Finnish FinBERT [Virtanen et al., 2019], Hebrew AlephBERT [Seker et al.,

2021], Hungarian huBERT [Nemeskey, 2021], Korean bert-kor-base, Persian ParsBERT [Farahani et al., 2021], Russian RuBert [Kuratov and Arkhipov, 2019], Swedish bert-base-swedish-cased, Turkish BERTurk, and Chinese bert-base-chinese models.

We only performed the experiments for a subset of studied languages for which we were able to find more specific BERT models. The aim is to check if the additional morphological features improve the performance of less general, i.e. more language-specific BERT models. These are trained on a lower number of languages, and larger amounts of texts in the included languages, compared to the original multilingual BERT model [Devlin et al., 2019] that was trained on 104 languages simultaneously. Due to this language-specific training, we expect these BERT models to capture the nuances of the languages better, thus possibly benefiting less from the additional morphological features.

We show the results in Table 9. In most cases, the specific multilingual BERT models, even without additional features, do as well as or better than the best performing original multilingual BERT model with additional features. The only worse LAS scores are achieved on Russian. This indicates that the more specific multilingual BERT models are generally better suited for the DP task than the original multilingual models. The addition of morphological features increases the LAS even further. The improvements range from 0.54% (Slovene) to 2.23% (Estonian) and are statistically significant for four out of six languages.

Table 9: LAS scores in the DP task achieved by more specific multilingual (top), and monolingual (bottom) BERT models. The more specific multilingual BERT models were trained on a smaller set of languages (three or five) than the original multilingual BERT model (104). For the base parsers with BERT embeddings, we display the improvement in LAS over using the original multilingual BERT for 104 languages (\uparrow_{mtl}). Statistically significant differences of best scores to mBERT baselines at $p = 0.01$ level are underlined.

lang.	model handle	+mBERT	(\uparrow_{mtl})	+mBERT +UPOS	+mBERT +UPOS +feats	+mBERT +feats
HRV	CroSloEngual BERT	87.84	(+1.47)	<u>89.01</u>	88.33	87.98
ENG	CroSloEngual BERT	88.37	(+0.28)	<u>89.69</u>	89.65	88.58
SLV	CroSloEngual BERT	93.98	(+2.31)	94.27	<u>94.52</u>	94.05
FIN	FinEst BERT	89.35	(+3.28)	<u>90.66</u>	90.61	90.54
EST	FinEst BERT	86.51	(+2.21)	88.67	<u>88.74</u>	87.44
RUS	Slavic BERT	85.60	(-0.39)	86.32	<u>86.39</u>	86.39
ARA	arabic-bert-base	83.40	(+2.43)	84.08	<u>84.93</u>	84.15
ENG	bert-base-cased	88.61	(+0.52)	<u>89.69</u>	89.56	88.59
EST	EstBERT	86.67	(+2.37)	88.69	<u>88.80</u>	87.32
FIN	FinBERT	91.51	(+5.44)	<u>92.19</u>	91.73	91.80
HEB	AlephBERT	89.41	(+2.23)	<u>90.37</u>	89.31	89.71
HUN	huBERT	81.28	(+4.21)	<u>82.88</u>	81.93	81.83
KOR	bert-kor-base	88.51	(+3.05)	<u>89.92</u>	/	/
FAS	ParsBERT	91.10	(+2.04)	<u>92.22</u>	92.20	91.69
RUS	RuBert	86.39	(+0.40)	<u>87.31</u>	87.04	86.61
SWE	bert-base-swedish-cased	89.41	(+2.74)	<u>90.63</u>	90.53	90.21
TUR	BERTurk	76.78	(+6.18)	76.75	<u>76.87</u>	76.07
ZHO	bert-base-chinese	84.72	(+5.23)	<u>86.31</u>	85.81	84.17

The monolingual models without additional features set an even higher baseline performance. The results of including additional information are mixed, though surprisingly many languages still see a significant increase in LAS. Out of twelve languages, the improvements are significant for eight and not significant for four languages.

These results indicate that the additional morphological features contain valuable information for the DP task, which the more specific BERT models still do not capture entirely. However, they improve over massively multilingual (i.e. more general) models. We suspect the increase would be even less pronounced if we experimented with “large“ (as opposed to base-sized) variants, e.g., the large English BERT would likely benefit even less from additional morphological information. We leave this line of experiments for further work.

6 Conclusion

We analysed adding explicit morphological information in the form of embeddings for POS tags and universal features to two currently dominant neural network architectures used in NLP: LSTM networks and transformer-based BERT models. We compared models enhanced with morphological information with baselines on three tasks (NER, DP, and CF). To obtain general conclusions, we used a variety of morphologically-rich languages from different language families. We make the code to re-run our experiments publicly available³.

The results indicate that adding morphological information to CF prediction models is not beneficial, but it improves the performance in the NER and DP tasks. For the DP task, the improvement depends on the quality of the morphological features. The additional morphological features consistently benefited LSTM-based models for NER and DP, both when they were of high quality and predicted (noisy). For BERT-based models, the *predicted* features do not make any practical difference for the NER and DP task but improve the performance in the DP task when they are of *high quality*. Testing different variants of BERT shows that language-specialised variants enhance the performance on the DP task and the additional morphological information is still beneficial, although less and less as we shift from multilingual towards monolingual models.

Comparing different BERT variants indicates that BERT models do not entirely capture the language morphology. Since the release of BERT, several new pre-training objectives have been proposed, such as syntactic and semantic phrase masking [Zhou et al., 2020b] and span masking [Joshi et al., 2020]. In further work, it makes sense to apply these models to the DP task to test how well they capture the morphology. Further, the effect of morphological features could be analysed on additional tasks and languages since the explicit morphological information does not seem to benefit them equally.

Acknowledgements

This work was supported by European Union’s Horizon 2020 Programme project EMBEDDIA (Cross-Lingual Embeddings for Less-Represented Languages in European News Media, grant no. 825153). The research was supported by the Slovene Research Agency through research core funding no. P6-0411 and the young researcher grant as well the Ministry of Culture of the Republic of Slovenia through project Development of Slovene in Digital Environment (RSDO). The Titan X Pascal used for a part of this research was donated by the NVIDIA Corporation.

References

- Mark Anderson and Carlos Gómez-Rodríguez. On the frailty of universal POS tags for neural UD parsers. In *Proceedings of the 24th Conference on Computational Natural Language Learning*, pages 69–96, 2020. doi: 10.18653/v1/2020.conll-1.6.
- 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, 2019. doi: 10.18653/v1/W19-3712.
- Miguel Ballesteros, Chris Dyer, and Noah A. Smith. Improved transition-based parsing by modeling characters instead of words with LSTMs. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 349–359, Lisbon, Portugal, September 2015. Association for Computational Linguistics. doi: 10.18653/v1/D15-1041.
- Yassine Benajiba, Paolo Rosso, and José Miguel Benedí Ruiz. ANERsys: An Arabic named entity recognition system based on maximum entropy. In Alexander Gelbukh, editor, *Computational Linguistics and Intelligent Text Processing*, pages 143–153, 2007.
- 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.
- Danqi Chen and Christopher Manning. A fast and accurate dependency parser using neural networks. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 740–750, 2014.
- Ronan Collobert and Jason Weston. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th International Conference on Machine Learning*, page 160–167, 2008.

³<https://github.com/matejklemen/morphological-dependency-parsing> (DP),
<https://github.com/EMBEDDIA/morphological-fasttext> (NER),
<https://github.com/EMBEDDIA/morphological-BERT> (NER),
<https://github.com/matejklemen/morphological-comment-filtering> (CF)

- Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. What you can cram into a single \$&!#* vector: Probing sentence embeddings for linguistic properties. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2126–2136, 2018.
- 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, 2019.
- Cícero dos Santos and Victor Guimarães. Boosting named entity recognition with neural character embeddings. In *Proceedings of the Fifth Named Entity Workshop*, pages 25–33, 2015. doi: 10.18653/v1/W15-3904.
- Timothy Dozat and Christopher D. Manning. Deep biaffine attention for neural dependency parsing. In *Proceedings on International Conference on Learning Representation*, 2016.
- Timothy Dozat, Peng Qi, and Christopher D. Manning. Stanford’s graph-based neural dependency parser at the CoNLL 2017 shared task. In *Proceedings of the CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 20–30, 2017. doi: 10.18653/v1/K17-3002.
- Daniel Edmiston. A systematic analysis of morphological content in BERT models for multiple languages. arXiv:2004.03032, 2020.
- Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav Goldberg. Amnesic probing: Behavioral explanation with amnesic counterfactuals. *Transactions of the Association for Computational Linguistics*, 9:160–175, 2021.
- Bojan Evkoski, Igor Mozetič, Nikola Ljubešić, and Petra Kralj Novak. Community evolution in retweet networks. *PLOS ONE*, 16(9):1–21, 09 2021. doi: 10.1371/journal.pone.0256175.
- Mehrdad Farahani, Mohammad Gharachorloo, Marzieh Farahani, and Mohammad Manthouri. ParsBERT: Transformer-based model for Persian language understanding. *Neural Processing Letters*, pages 1–17, 2021.
- Paula Fortuna and Sérgio Nunes. A survey on automatic detection of hate speech in text. *ACM Computing Surveys*, 51(4), 2018.
- Lei Gao and Ruihong Huang. Detecting online hate speech using context aware models. In *Proceedings of the International Conference Recent Advances in Natural Language Processing, RANLP 2017*, pages 260–266, 2017.
- Spiros V Georgakopoulos, Sotiris K Tasoulis, Aristidis G Vrahatis, and Vassilis P Plagianakos. Convolutional neural networks for toxic comment classification. In *Proceedings of the 10th Hellenic Conference on Artificial Intelligence*, pages 1–6, 2018.
- Stefan Grünewald, Annemarie Friedrich, and Jonas Kuhn. Applying Occam’s razor to transformer-based dependency parsing: What works, what doesn’t, and what is really necessary. In *Proceedings of the 17th International Conference on Parsing Technologies and the IWPT 2021 Shared Task on Parsing into Enhanced Universal Dependencies (IWPT 2021)*, pages 131–144, 2021. doi: 10.18653/v1/2021.iwpt-1.13.
- Onur Güngör, Eray Yıldız, Suzan Üsküdarlı, and Tunga Güngör. Morphological embeddings for named entity recognition in morphologically rich languages. *arXiv preprint arXiv:1706.00506*, 2017.
- Onur Güngör, Tunga Güngör, and Suzan Üsküdarlı. The effect of morphology in named entity recognition with sequence tagging. *Natural Language Engineering*, 25(1):147–169, 2019. doi: 10.1017/S1351324918000281.
- Jan Hajič and Dan Zeman, editors. *Proceedings of the CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*. Association for Computational Linguistics, 2017.
- Jiahui Han, Shengtian Wu, and Xinyu Liu. jhan014 at SemEval-2019 task 6: Identifying and categorizing offensive language in social media. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 652–656, June 2019.
- Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9(8):1735–1780, 1997.
- Zhiheng Huang, Wei Xu, and Kai Yu. Bidirectional LSTM-CRF models for sequence tagging. *ArXiv*, abs/1508.01991, 2015.
- Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. What does BERT learn about the structure of language? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3651–3657, 2019.
- Tao Ji, Yuanbin Wu, and Man Lan. Graph-based dependency parsing with graph neural networks. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2475–2485, July 2019.
- Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. SpanBERT: Improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics*, 8: 64–77, 2020.

- Daniel Jurafsky and James H. Martin. *Speech and Language Processing (2nd Edition)*. Prentice-Hall, Inc., USA, 2009.
- Gopal K Kanji. *100 statistical tests*. Sage, 2006.
- Jurgita Kapočiušė-Dzikienė, Joakim Nivre, and Algis Krupavičius. Lithuanian dependency parsing with rich morphological features. In *Proceedings of the fourth workshop on statistical parsing of morphologically-rich languages*, pages 12–21, 2013.
- Mojtaba Khallash, Ali Hadian, and Behrouz Minaei-Bidgoli. An empirical study on the effect of morphological and lexical features in Persian dependency parsing. In *Proceedings of the Fourth Workshop on Statistical Parsing of Morphologically-Rich Languages*, pages 97–107, 2013.
- Eliyahu Kiperwasser and Yoav Goldberg. Simple and accurate dependency parsing using bidirectional LSTM feature representations. *Transactions of the Association for Computational Linguistics*, 4:313–327, 2016.
- Dan Kondratyuk and Milan Straka. 75 languages, 1 model: Parsing universal dependencies universally. 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 2779–2795, 2019.
- Simon Krek, Kaja Dobrovoljc, Tomaž Erjavec, Sara Može, Nina Ledinek, Nanika Holz, Katja Zupan, Polona Gantar, Taja Kuzman, Jaka Čibej, Špela Arhar Holdt, Teja Kavčič, Iza Škrjanec, Dafne Marko, Lucija Jezeršek, and Anja Zajc. Training corpus ssj500k 2.2, 2019. URL <http://hdl.handle.net/11356/1210>. Slovenian language resource repository CLARIN.SI.
- Artur Kulmizev, Miryam de Lhoneux, Johannes Gontrum, Elena Fano, and Joakim Nivre. Deep contextualized word embeddings in transition-based and graph-based dependency parsing - a tale of two parsers revisited. 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 2755–2768, 2019.
- Yuri Kuratov and Mikhail Arkipov. Adaptation of deep bidirectional multilingual transformers for Russian language. In *Computational Linguistics and Intellectual Technologies: Proceedings of the International Conference “Dialogue 2019”*, 2019.
- Onur Kuru, Ozan Arkan Can, and Deniz Yuret. CharNER: Character-level named entity recognition. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 911–921, 2016.
- Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. Neural architectures for named entity recognition. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 260–270, 2016. doi: 10.18653/v1/N16-1030.
- Gina-Anne Levow. The third international Chinese language processing bakeoff: Word segmentation and named entity recognition. In *Proceedings of the Fifth SIGHAN Workshop on Chinese Language Processing*, pages 108–117, 2006.
- Miryam de Lhoneux, Yan Shao, Ali Basirat, Eliyahu Kiperwasser, Sara Stymne, Yoav Goldberg, and Joakim Nivre. From raw text to Universal Dependencies - look, no tags! In *Proceedings of the CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 207–217, 2017. doi: 10.18653/v1/K17-3022.
- Zuchao Li, Jiaxun Cai, Shexia He, and Hai Zhao. Seq2seq dependency parsing. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 3203–3214, 2018.
- KyungTae Lim, Cheoneum Park, Changki Lee, and Thierry Poibeau. SEx BiST: A multi-source trainable parser with deep contextualized lexical representations. In *Proceedings of the CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 143–152, 2018. doi: 10.18653/v1/K18-2014.
- KyungTae Lim, Jay Yoon Lee, Jaime Carbonell, and Thierry Poibeau. Semi-supervised learning on meta structure: Multi-task tagging and parsing in low-resource scenarios. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8344–8351, 2020. doi: 10.1609/aaai.v34i05.6351.
- Yongjie Lin Lin, Yi Chern Tan, and Robert Frank. Open Sesame: Getting inside BERT’s linguistic knowledge. In *Proceedings of the Second BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, 2019.
- Nikola Ljubešić, Željko Agić, Filip Klubička, Vuk Batanović, and Tomaž Erjavec. Training corpus hr500k 1.0, 2018. URL <http://hdl.handle.net/11356/1183>. Slovenian language resource repository CLARIN.SI.
- Shervin Malmasi and Marcos Zampieri. Detecting hate speech in social media. In *Proceedings of the International Conference Recent Advances in Natural Language Processing*, pages 467–472, 2017.
- Yuval Marton, Nizar Habash, and Owen Rambow. Improving Arabic dependency parsing with lexical and inflectional morphological features. In *Proceedings of the NAACL HLT 2010 First Workshop on Statistical Parsing of Morphologically-Rich Languages*, pages 13–21, 2010.

- Ryan McDonald, Fernando Pereira, Kiril Ribarov, and Jan Hajič. Non-projective dependency parsing using spanning tree algorithms. In *Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing*, pages 523–530, 2005.
- Vladislav Mikhailov, Oleg Serikov, and Ekaterina Artemova. Morph call: Probing morphosyntactic content of multilingual transformers. In *Proceedings of the Third Workshop on Computational Typology and Multilingual NLP*, pages 97–121, 2021. doi: 10.18653/v1/2021.sigtyp-1.10.
- Kristian Miok, Dong Nguyen-Doan, Blaž Škrlj, Daniela Zaharie, and Marko Robnik-Šikonja. Prediction uncertainty estimation for hate speech classification. In *International Conference on Statistical Language and Speech Processing*, pages 286–298, 2019.
- Mahdi Mohseni and Amirhossein Tebbifakhr. MorphoBERT: A Persian NER system with BERT and morphological analysis. In *Proceedings of The First International Workshop on NLP Solutions for Under Resourced Languages (NSURL 2019) co-located with ICNLSP 2019 - Short Papers*, pages 23–30, 11–12 September 2019.
- Jihyung Moon, Won Ik Cho, and Junbum Lee. BEEP! Korean corpus of online news comments for toxic speech detection. In *Proceedings of the Eighth International Workshop on Natural Language Processing for Social Media*, pages 25–31, 2020.
- Dávid Márk Nemeskey. Introducing huBERT. In *XVII. Magyar Számítógépes Nyelvészeti Konferencia (MSZNY2021)*, 2021.
- Dat Quoc Nguyen and Karin Verspoor. An improved neural network model for joint POS tagging and dependency parsing. In *Proceedings of the CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 81–91, Brussels, Belgium, October 2018. Association for Computational Linguistics. doi: 10.18653/v1/K18-2008.
- Linh The Nguyen and Dat Quoc Nguyen. PhoNLP: A joint multi-task learning model for Vietnamese part-of-speech tagging, named entity recognition and dependency parsing. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Demonstrations*, pages 1–7, 2021. doi: 10.18653/v1/2021.naacl-demos.1.
- Joakim Nivre. An efficient algorithm for projective dependency parsing. In *Proceedings of the Eighth International Conference on Parsing Technologies*, pages 149–160, 2003.
- Joakim Nivre, Mitchell Abrams, Željko Agić, et al. Universal Dependencies 2.6, 2020. URL <http://hdl.handle.net/11234/1-2988>. LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.
- Şaziye Betül Özateş, Arzucan Özgür, Tunga Güngör, and Balkız Öztürk. A morphology-based representation model for LSTM-based dependency parsing of agglutinative languages. In *Proceedings of the CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 238–247, Brussels, Belgium, October 2018. doi: 10.18653/v1/K18-2024.
- Pēteris Paikens, Ilze Auziņa, Ginta Garkaje, and M Paegle. Towards named entity annotation of Latvian national library corpus. *Frontiers in Artificial Intelligence and Applications*, 247:169–175, 01 2012. doi: 10.3233/978-1-61499-133-5-169.
- Wenzhe Pei, Tao Ge, and Baobao Chang. An effective neural network model for graph-based dependency parsing. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 313–322, 2015.
- Telmo Pires, Eva Schlinger, and Dan Garrette. How multilingual is multilingual BERT? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4996–5001, 2019.
- Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. Stanza: A Python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, 2020.
- Anna Rogers, Olga Kovaleva, and Anna Rumshisky. A primer in BERTology: What we know about how BERT works. *Transactions of the Association for Computational Linguistics*, 8:842–866, 2020.
- Teemu Ruokolainen, Pekka Kauppinen, Miikka Silfverberg, and Krister Lindén. A Finnish news corpus for named entity recognition. *Language Resources and Evaluation*, 54, 08 2019. doi: 10.1007/s10579-019-09471-7.
- Ali Safaya, Moutasem Abdullatif, and Deniz Yuret. KUISAIL at SemEval-2020 task 12: BERT-CNN for offensive speech identification in social media. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 2054–2059, 2020.

- Tatjana Scheffler, Erik Haegert, Santichai Pornavalai, and Mino Lee Sasse. Feature explorations for hate speech classification. In *14th Conference on Natural Language Processing KONVENS 2018*, volume 6, page 8, 2018.
- Djame Seddah, Sandra Koebler, and Reut Tsarfaty, editors. *Proceedings of the NAACL HLT 2010 First Workshop on Statistical Parsing of Morphologically-Rich Languages*, 2010.
- Wolfgang Seeker and Jonas Kuhn. On the role of explicit morphological feature representation in syntactic dependency parsing for German. In *Proceedings of the 12th International Conference on Parsing Technologies*, pages 58–62, 2011.
- Amit Seker, Elron Bandel, Dan Bareket, Idan Brusilovsky, Refael Shaked Greenfeld, and Reut Tsarfaty. AlephBERT: A Hebrew large pre-trained language model to start-off your Hebrew NLP application with. ArXiv 2104.04052, 2021.
- Serhiy Shtovba, Olena Shtovba, and Mykola Petrychko. Detection of social network toxic comments with usage of syntactic dependencies in the sentences. In *Proceedings of the Second International Workshop on Computer Modeling and Intelligent Systems*, pages 313–323, 2019.
- Lilia Simeonova, Kiril Simov, Petya Osenova, and Preslav Nakov. A morpho-syntactically informed LSTM-CRF model for named entity recognition. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2019)*, pages 1104–1113, 2019.
- Anatoli Starostin, Victor V. Bocharov, Svetlana Alexeeva, A. A. Bodrova, Alexander Chuchunkov, Sh. Sh. Dzhumaev, Irina Efimenko, D V Granovsky, Vladimir F. Khoroshevsky, Irina V. Krylova, Marina Nikolaeva, Ivan Smurov, and Svetlana Toldova. FactRuEval 2016: Evaluation of named entity recognition and fact extraction systems for Russian. In *Annual International Conference “Dialogue“*, 2016.
- Milan Straka. UDPipe 2.0 prototype at CoNLL 2018 UD shared task. In *Proceedings of the CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 197–207, 2018. doi: 10.18653/v1/K18-2020.
- Jana Straková, Milan Straka, and Jan Hajič. Neural networks for featureless named entity recognition in Czech. In *Proceedings of Text, Speech, and Dialogue*, pages 173–181, 2016.
- Student. The probable error of a mean. *Biometrika*, pages 1–25, 1908.
- Nasrin Taghizadeh, Zeinab Borhanifard, Melika Golestani Pour, Mojgan Farhoodi, Maryam Mahmoudi, Masoumeh Azimzadeh, and Hesham Faili. NSURL-2019 task 7: Named entity recognition for Farsi. In *Proceedings of The First International Workshop on NLP Solutions for Under Resourced Languages (NSURL 2019) co-located with ICNLP 2019 - Short Papers*, pages 9–15, 2019.
- Hasan Tanvir, Claudia Kittask, Sandra Eiche, and Kairit Sirts. EstBERT: A pretrained language-specific BERT for Estonian. In *Proceedings of the 23rd Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 11–19, 2021.
- Ian Tenney, Dipanjan Das, and Ellie Pavlick. BERT rediscovers the classical NLP pipeline. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601, 2019.
- Erik F. Tjong Kim Sang and Fien De Meulder. Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition. In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*, pages 142–147, 2003.
- Alexander Tkachenko, Timo Petmanson, and Sven Laur. Named entity recognition in Estonian. In *Proceedings of the 4th Biennial International Workshop on Balto-Slavic Natural Language Processing*, pages 78–83, 2013.
- Matej Ulčar and Marko Robnik-Šikonja. FinEst BERT and CroSloEngual BERT: less is more in multilingual models. In *Proceedings of Text, Speech, and Dialogue TSD 2020*, 2020. doi: 10.1007/978-3-030-58323-1_11.
- Cynthia Van Hee, Els Lefever, Ben Verhoeven, Julie Mennes, Bart Desmet, Guy De Pauw, Walter Daelemans, and Véronique Hoste. Detection and Fine-Grained Classification of Cyberbullying Events. In *Proceedings of the International Conference Recent Advances in Natural Language Processing*, pages 672–680, 2015.
- Clara Vania, Andreas Grivas, and Adam Lopez. What do character-level models learn about morphology? The case of dependency parsing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2573–2583, 2018. doi: 10.18653/v1/D18-1278.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, page 6000–6010, 2017.
- 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 1912.07076, 2019.

- Alexsander Waibel, Toshiyuki Hanazawa, Geoffrey Hinton, Kiyohiro Shikano, and Kevin J. Lang. Phoneme recognition using time-delay neural networks. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 37(3):328–339, 1989. doi: 10.1109/29.21701.
- Frank Wilcoxon, SK Katti, and Roberta A Wilcox. Critical values and probability levels for the Wilcoxon rank sum test and the Wilcoxon signed rank test. *Selected tables in mathematical statistics*, 1:171–259, 1970.
- Ellery Wulczyn, Nithum Thain, and Lucas Dixon. Ex machina: Personal attacks seen at scale. In *Proceedings of the 26th International Conference on World Wide Web*, page 1391–1399, 2017. doi: 10.1145/3038912.3052591.
- Hiroyasu Yamada and Yuji Matsumoto. Statistical dependency analysis with support vector machines. In *Proceedings of the Eighth International Conference on Parsing Technologies*, pages 195–206, 2003.
- Zhilin Yang, Ruslan Salakhutdinov, and William Cohen. Multi-task cross-lingual sequence tagging from scratch, 2016.
- Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. SemEval-2019 task 6: Identifying and categorizing offensive language in social media (OffensEval). In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 75–86, 2019.
- Marcos Zampieri, Preslav Nakov, Sara Rosenthal, Pepa Atanasova, Georgi Karadzhov, Hamdy Mubarak, Leon Derczynski, Zeses Pitenis, and Çağrı Çöltekin. SemEval-2020 Task 12: Multilingual Offensive Language Identification in Social Media (OffensEval 2020). In *Proceedings of SemEval*, 2020.
- Houquan Zhou, Yu Zhang, Zhenghua Li, and Min Zhang. Is POS tagging necessary or even helpful for neural dependency parsing? In *Natural Language Processing and Chinese Computing - 9th CCF International Conference, NLPCC 2020, Proceedings, Part I*, pages 179–191, 2020a. doi: 10.1007/978-3-030-60450-9_15.
- Junru Zhou, Zhuosheng Zhang, Hai Zhao, and Shuailiang Zhang. LIMIT-BERT : Linguistics informed multi-task BERT. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4450–4461, 2020b. doi: 10.18653/v1/2020.findings-emnlp.399.

A Used Stanza models

(a) Stanza models used in the NER experiments. (b) Stanza models used in the additional DP experiments. (c) Stanza models used in the CF experiments.

lang.	Model	lang.	Model	accuracy (UPOS)	accuracy (UFeats)	lang.	Model
Arabic	padt	English	gum	92.45	93.92	Arabic	padt
Chinese	gsdsimp	Estonian	ewt	91.15	88.86	English	ewt
Croatian	set	Finnish	ftb	87.59	86.20	Greek	gdt
English	ewt	Korean	gsd	70.81	/	Korean	gsd
Estonian	edt	Persian	seraji	82.32	76.30	Slovene	ssj
Finnish	tdt	Russian	syntagrus	89.32	85.22	Turkish	imst
Korean	kaist	Slovene	sst	80.45	78.74		
Latvian	lvtb	Swedish	lines	94.66	86.16		
Russian	syntagrus	Turkish	imst	86.20	68.83		
Slovene	ssj						
Swedish	talbanken						

B Results of dependency parsing with cased multilingual BERT

Table 11 shows the results of a subset of the main dependency parsing experiments performed using a **cased** multilingual BERT model (bert-base-multilingual-cased). However, the conclusions drawn from Table 5 and Table 11 are the same.

Table 11: LAS achieved using a cased (instead of uncased) multilingual BERT model on a random sample of eight languages in the DP task. Statistically significant differences between baselines and best scores at $p = 0.01$ level are underlined.

lang.	+mBERT	+mBERT +UPOS	+mBERT +UPOS +feats	+mBERT +feats
ZHO	83.03	<u>84.97</u>	84.23	82.53
ENG	88.04	<u>88.74</u>	88.59	88.19
FIN	85.65	<u>87.55</u>	87.34	87.17
HUN	77.77	80.14	<u>80.40</u>	77.80
FAS	88.99	<u>90.66</u>	90.57	89.96
SLV	91.35	92.46	<u>93.16</u>	92.09
SWE	85.66	<u>87.23</u>	87.16	86.30
TUR	71.52	<u>72.82</u>	72.38	70.92