

# Sentiment Analysis with Contextual Embeddings and Self-Attention

Katarzyna Biesialska <sup>\*1</sup>[0000–0002–2865–7990],  
Magdalena Biesialska<sup>\*1</sup>[0000–0001–7890–3523], and  
Henryk Rybinski<sup>2</sup>[0000–0002–2890–7080]

<sup>1</sup> Universitat Politècnica de Catalunya, Barcelona, Spain  
{katarzyna,magdalena}.biesialska@upc.edu

<sup>2</sup> Warsaw University of Technology, Warsaw, Poland h.rybinski@ii.pw.edu.pl

**Abstract.** In natural language the intended meaning of a word or phrase is often implicit and depends on the context. In this work, we propose a simple yet effective method for sentiment analysis using contextual embeddings and a self-attention mechanism. The experimental results for three languages, including morphologically rich Polish and German, show that our model is comparable to or even outperforms state-of-the-art models. In all cases the superiority of models leveraging contextual embeddings is demonstrated. Finally, this work is intended as a step towards introducing a universal, multilingual sentiment classifier.

**Keywords:** Sentiment classification · Deep learning · Word embeddings.

## 1 Introduction

All areas of human life are affected by people’s views. With the sheer amount of reviews and other opinions over the Internet, there is a need for automating the process of extracting relevant information. For machines, however, measuring sentiment is not an easy task, because natural language is highly ambiguous at all levels, and thus difficult to process. For instance, a single word can hardly convey the whole meaning of a statement. Moreover, computers often do not distinguish literal from figurative meaning or incorrectly handle complex linguistic phenomena, such as: sarcasm, humor, negation etc.

In this paper, we take a closer look at two factors that make automatic opinion mining difficult – the problem of representing text information, and sentiment analysis (SA). In particular, we leverage contextual embeddings, which enable to convey a word meaning depending on the context it occurs in. Furthermore, we build a hierarchical multi-layer classifier model, based on an architecture of the Transformer encoder [32], primarily relying on a self-attention mechanism and bi-attention. The proposed sentiment classification model is language independent, which is especially useful for low-resource languages (e.g. Polish).

---

\* Both authors contributed equally to this work, which was mostly done at the Warsaw University of Technology.

We evaluate our methods on various standard datasets, which allows us to compare our approach against current state-of-the-art models for three languages: English, Polish and German. We show that our approach is comparable to the best performing sentiment classification models; and, importantly, in two cases yields significant improvements over the state of the art.

The paper is organized as follows: Section 2 presents the background and related work. Section 3 describes our proposed method. Section 4 discusses datasets, experimental setup, and results. Section 5 concludes this paper and outlines the future work.

## 2 Related Work

Sentiment classification has been one of the most active research areas in natural language processing (NLP) and has become one of the most popular downstream tasks to evaluate performance of neural network (NN) based models. The task itself encompasses several different opinion related tasks, hence it tackles many challenging NLP problems, see e.g. [16, 20].

### 2.1 Sentiment Analysis Approaches

The first fully-formed techniques for SA emerged around two decades ago, and continued to be prevalent for several years, until deep learning methods entered the stage. The most straight-forward method, developed in [30], is based on the number of positive and negative words in a piece of text. Concretely, the text is assumed to have positive polarity if it contains more positive than negative terms, and vice versa. Of course, the term-counting method is often insufficient; therefore, an improved method was proposed in [10], which combines counting positive and negative terms with a machine learning (ML) approach (i.e. Support Vector Machine).

Various studies (e.g. [31]) have shown that one can determine the polarity of an unknown word by calculating co-occurrence statistics of it. Moreover, classical solutions to the SA problem are often based on lexicons. Traditional lexicon-based SA leverages word-lists, that are pre-annotated with positive and negative sentiment. Therefore, for many years lexicon-based approaches have been utilized when there was insufficient amount of labeled data to train a classifier in a fully supervised way.

In general, ML algorithms are popular methods for determining sentiment polarity. A first ML model applied to SA has been implemented in [21]. Moreover, throughout the years, different variants of NN architectures have been introduced in the field of SA. Especially recursive neural networks [22], such as recurrent neural networks (RNN) [28, 29, 13], or convolutional neural networks (CNN) [9, 11] have become the most prevalent choices.

## 2.2 Vector Representations of Words

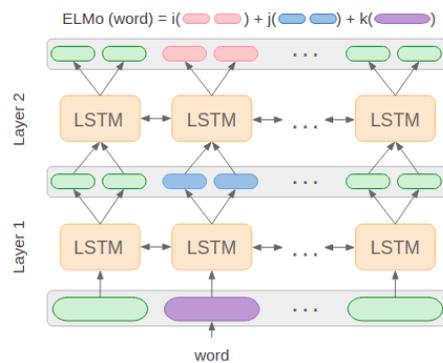
One of the principal concepts in linguistics states that related words can be used in similar ways [6]. Importantly, words may have different meaning in different contexts. Nevertheless, until recently it has been a dominant approach (e.g. word2vec [19], GloVe [23]) to learn representations such that each and every word has to capture all its possible meanings.

However, lately a new set of methods to learn dynamic representations of words has emerged [18, 7, 24, 25, 5]. These approaches allow each word representation to capture what a word means in a particular context. While every word token has its own vector, the vector can depend on a variable-length sequence of nearby words (i.e. context). Consequently, a context vector is obtained by feeding a neural network with these context word vectors and subsequently encoding them into a single fixed-length vector.

ULMFiT [7] was the very first method to induce contextual word representations by harnessing the power of language modeling. The authors proposed to learn contextual embeddings by pre-training a language model (LM), and then performing task-specific fine-tuning. ULMFiT architecture is based on a vanilla 3-layer Long Short-Term Memory (LSTM) NN without any attention mechanism.

The other contextual embedding model introduced recently is called ELMo (Embeddings from Language Models) [24]. Similarly to ULMFiT, this model uses tokens at the word-level. ELMo contextual embeddings are “deep” as they are a function of all hidden states. Concretely, context-sensitive features are extracted from a left-to-right and a right-to-left 2-layer bidirectional LSTM language models. Thus, the contextual representation of each word is the concatenation of the left-to-right and right-to-left representations as well as the initial embedding (see Fig. 1).

The most recent model – BERT [5] – is more sophisticated architecturally-wise, as it is a multi-layer masked LM based on the Transformer NN utilizing sub-word tokens. However, as we are bound to use word-level tokens in our sentiment classifier, we leverage the ELMo model for obtaining contextual embeddings. More specifically, by means of ELMo we are able to feed our classifier model with context-aware embeddings of an input sequence. Hence, in this setting we do not perform any fine-tuning of ELMo on a downstream task.



**Fig. 1.** The architecture of ELMo.

### 2.3 Self-Attention Deep Neural Networks

The attention mechanism was introduced in [3] in 2014 and since then it has been applied successfully to different computer vision (e.g. visual explanation) and NLP (e.g. machine translation) tasks. The mechanism is often used as an extra source of information added on top of the CNN or LSTM model to enhance the extraction of sentence embedding [26, 15]. However, this scenario is not applicable to sentiment classification, since the model only receives a single sentence on input, hence there is no such extra information [15].

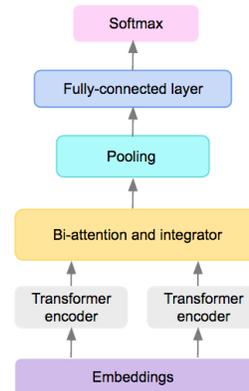
Self-attention (or intra-attention) is an attention mechanism that computes a representation of a sequence by relating different positions of a single sequence. Previous work on sentiment classification has not covered extensively attention-based neural network models for SA (especially using the Transformer architecture [32]), although some papers have appeared recently [2, 14].

## 3 The Proposed Approach

Our proposed model, called Transformer-based Sentiment Analysis (TSA) (see Fig. 2), is based on the recently introduced Transformer architecture [32], which has provided significant improvements for the neural machine translation task. Unlike RNN or CNN based models, the Transformer is able to learn dependencies between distant positions. Therefore, in this paper we show that attention-based models are suitable for other NLP tasks, such as learning distributed representations and sentiment analysis, and thus are able to improve the overall accuracy.

The architecture of the TSA model and steps to train it can be summarized as follows:

- a) At the very beginning there is a simple text pre-processing method that performs text clean-up and splits text into tokens.
- b) We use contextual word representations to represent text as real-valued vectors.
- c) After embedding the text into real-valued vectors, the Transformer network maps the input sequence into hidden states using self-attention.
- d) Next a bi-attention mechanism is utilized to estimate the interdependency between representations.
- e) A single layer LSTM together with self-attentive pooling compute the pooled representations.
- f) A joint representation for the inputs is later passed to a fully-connected neural network.
- g) Finally, a softmax layer is used to determine sentiment of the text.



**Fig. 2.** An overview of the TSA model architecture.

### 3.1 Embeddings and Encoded Positional Information

Non-recurrent models, such as deep self-attention NN, do not necessarily process the input sequence in a sequential manner. Hence, there is no way they can record the position of each word in a sequence, which is an inherent limitation of every such model. Therefore, in the case of the Transformer, the need has been addressed in the following manner – the Transformer takes into account the order of the words in the input sequence by encoding their position information in extra vectors (so called positional encoding vectors) and adding them to input embeddings. There are many different approaches to embed position information, such as learned or fixed positional encodings (PE), or recently introduced relative position representations (RPR) [27]. The original Transformer used sine and cosine functions of different frequencies.

In this work, we explore the effectiveness of applying a modified approach to incorporate positional information into the model, namely using RPR instead of PE. Furthermore, we use global average pooling in order to average the output of the last self-attention layer and prepare the model for the final classification layer.

### 3.2 The Transformer Encoder

The input sequence is combined with word and positional embeddings, which provide time signal, and together are fed into an encoder block. Matrices for a query  $Q$ , a key  $K$  and a value  $V$  are calculated and passed to a self-attention layer. Next, a normalization is applied and residual connections provide additional context. Further, a final dense layer with vocabulary size generates the output of the encoder. A fully-connected feed-forward network within the model is a single hidden layer network with a ReLU activation function in between:

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2 \quad (1)$$

### 3.3 Self-Attention Layer

The self-attention block in the encoder is called multi-head self-attention. A self-attention layer allows each position in the encoder to access all positions in the previous layer of the encoder immediately, and in the first layer all positions in the input sequence. The multi-head self-attention layer employs  $h$  parallel self-attention layers, called heads, with different  $Q$ ,  $K$ ,  $V$  matrices obtained for each head. In a nutshell, the attention mechanism in the Transformer architecture relies on a scaled dot-product attention, which is a function of  $Q$  and a set of  $K$ - $V$  pairs. The computation of attention is performed in the following order. First, a multiplication of a query and transposed key is scaled through the scaling factor of  $1/\sqrt{d_z}$  (Eq. 2)

$$m_{ij} = \frac{QK^T}{\sqrt{d_z}} \quad (2)$$

Next, the attention is produced using the softmax function over their scaled inner product:

$$\alpha_{ij} = \frac{e^{m_{ij}}}{\sum_{k=1}^n e^{m_{ik}}} \quad (3)$$

Finally, the weighted sum of each attention head and a value is calculated as follows:

$$z_i = \sum_{j=1}^n \alpha_{ij} V \quad (4)$$

### 3.4 Masking and Pooling

Similar to other sources of data, the datasets used for training and evaluation of our models contain sequences of different length. The most common approach in the literature involves finding a maximal sequence length existing in the dataset/batch and padding sentences that are shorter than the longest one with trailing zeroes. In the proposed TSA model, we deal with the problem of variable-length sequences by using masking and self-attentive pooling. The inspiration for our approach comes from the BCN model proposed in [18]. Thanks to this mechanism, we are able to fit sequences of different length into the final fixed-size vector, which is required for the computation of the sentiment score. The self-attentive pooling layer is applied just after the encoder block.

## 4 Experiments

### 4.1 Datasets

In this work, we compare sentiment analysis results considering four benchmark datasets in three languages. All datasets are originally split into training, dev and test sets. Below we describe these datasets in more detail.

**Table 1.** Sentiment analysis datasets with number of classes and train/dev/test split.

Dataset	# Classes	Train	Dev	Test	Domain	Language
SST-2	2	6,920	872	1,821	movies	English
SST-5	5	8,544	1,101	2,210	movies	English
PolEmo 2.0-IN	5	5,783	723	722	medical, hotels	Polish
GermEval	3	19,432	2,369	2,566	travel, transport	German

**Stanford Sentiment Treebank (SST)** This collection of movie reviews [28] from the `rottentomatoes.com` is annotated for the binary (SST-2) and fine-grained (SST-5) sentiment classification. SST-2 divides reviews into two groups: *positive* and *negative*, while SST-5 distinguishes 5 different review types: *very positive*, *positive*, *neutral*, *negative*, *very negative*. The dataset consists of 11,855 single sentences and is widely used in the NLP community.

**PolEmo 2.0** The dataset [12] comprises online reviews from education, medicine and hotel domains. There are two separate test sets, to allow for in-domain (medicine and hotels) and out-of-domain (products and university) evaluation. The dataset comes with the following sentiment labels: *strong positive*, *weak positive*, *neutral*, *weak negative*, *strong negative*, and *ambiguous*.

**GermEval** This dataset [33] contains customer reviews of the railway operator (Deutsche Bahn) published on social media and various web pages. Customers expressed their feedback regarding the service of the railway company (e.g. travel experience, timetables, etc.) by rating it as *positive*, *negative*, or *neutral*.

## 4.2 Experimental Setup

Pre-processing of input datasets is kept to a minimum as we perform only tokenization when required. Furthermore, even though some datasets, such as SST or GermEval, provide additional information (i.e. phrase, word or aspect-level annotations), for each review we only extract text of the review and its corresponding rating.

The model is implemented in the Python programming language, PyTorch<sup>3</sup> and AllenNLP<sup>4</sup>. Moreover, we use pre-trained word-embeddings, such as ELMo [24], GloVe [23]. Specifically, we use the following ELMo models: Original<sup>5</sup>, Polish [8] and German [17]. In the ELMo+GloVe+BCN model we use the following 300-dimension GloVe embeddings: English<sup>6</sup>, Polish [4] and German<sup>7</sup>. In order to simplify our approach when training the sentiment classifier model, we establish a very similar setting to the vanilla Transformer. We use the same optimizer - Adam with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.98$ , and  $\epsilon = 10^{-9}$ . We incorporate four types of regularization during training: dropout probability  $P_{drop} = 0.1$ , embedding dropout probability  $P_{emb} = 0.5$ , residual dropout probability  $P_{res} = 0.2$ , and attention dropout probability  $P_{attn} = 0.1$ . We use 2 encoder layers. In addition, we employ label smoothing of value  $\epsilon_{ls} = 0.1$ . In terms of RPR parameters, we set clipping distance to  $k = 10$ .

## 4.3 Results and Discussion

In Table 2, we summarize experimental results achieved by our model and other state-of-the-art systems reported in the literature by their respective authors.

We observe that our models, baseline and ELMo+TSA, outperform state-of-the-art systems for all three languages. More importantly, the presented accuracy scores indicate that the TSA model is competitive and for two languages (Polish and German) achieves the best results. Also noteworthy, in Table 2,

<sup>3</sup> <https://pytorch.org>

<sup>4</sup> <https://allennlp.org>

<sup>5</sup> <https://allennlp.org/elmo>

<sup>6</sup> <http://nlp.stanford.edu/data/glove.840B.300d.zip>

<sup>7</sup> <https://wikipedia2vec.github.io/wikipedia2vec/pretrained>

**Table 2.** Results of our systems compared to baselines and state-of-the-art systems evaluated on English, Polish and German sentiment classification datasets.

	<b>English</b>		<b>Polish</b>	<b>German</b>
	SST-2	SST-5	PolEmo2.0-IN	GermEval
RNTN [28]	85.4	45.7	-	-
DCNN [9]	86.8	48.5	-	-
CNN [11]	88.1	48.0	-	-
DMN [13]	88.6	52.1	-	-
Constituency Tree-LSTM [29]	88.0	51.0	-	-
CoVe+BCN [18]	90.3	<b>53.7</b>	-	-
SSAN+RPR [2]	84.2	48.1	-	-
Polish BERT [1]	-	-	88.1	-
SWN2-RNN [33]	-	-	-	74.9
<i>Our baseline</i>				
ELMo+GloVe+BCN	<b>91.4</b>	53.5	88.9	78.2
<i>Our model</i>				
ELMo+TSA	89.3	50.6	<b>89.8</b>	<b>78.9</b>

there are two models that use some variant of the Transformer: SSAN+RPR [2] uses the Transformer encoder for the classifier, while Polish BERT [1] employs Transformer-based language model introduced in [5]. One of the reasons why we achieve higher score for the SST dataset might be that the authors of SSAN+RPR used word2vec embeddings [19], whereas we employ ELMo contextual embeddings [24]. Moreover, in our TSA model we use not only self-attention (as in SSAN+RPR) but also a bi-attention mechanism, hence this also should provide performance gains over standard architectures.

In conclusion, comparing the results of the models leveraging contextual embeddings (CoVe+BCN, Polish BERT, ELMo+GloVe+BCN and ELMo+TSA) with the rest of the reported models, which use traditional distributional word vectors, we note that the former category of sentiment classification systems demonstrates remarkably better results.

## 5 Conclusion and Future Work

We have presented a novel architecture, based on the Transformer encoder with relative position representations. Unlike existing models, this work proposes a model relying solely on a self-attention mechanism and bi-attention. We show that our sentiment classifier model achieves very good results, comparable to the state of the art, even though it is language-agnostic. Hence, this work is a step towards building a universal, multi-lingual sentiment classifier.

In the future, we plan to evaluate our model using benchmarks also for other languages. It is particularly interesting to analyze the behavior of our model

with respect to low-resource languages. Finally, other promising research avenues worth exploring are related to unsupervised cross-lingual sentiment analysis.

## References

1. Allegro: Klej benchmark, <https://klejbenchmark.com/>, accessed: 2020-01-20
2. Ambartsoumian, A., Popowich, F.: Self-attention: A better building block for sentiment analysis neural network classifiers. In: Proceedings of the 9th EMNLP Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis. pp. 130–139 (2018)
3. Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. arXiv (2014)
4. Dadas, S.: A repository of polish NLP resources. Github (2019), <https://github.com/sdadas/polish-nlp-resources/>, accessed: 2020-01-20
5. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: 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. pp. 4171–4186 (2019)
6. Firth, J.R.: A synopsis of linguistic theory, 1930-1955. Studies in linguistic analysis (1957)
7. Howard, J., Ruder, S.: Universal language model fine-tuning for text classification. In: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics. pp. 328–339 (2018)
8. Janz, A., Milkowski, P.: ELMo embeddings for polish (2019), <http://hdl.handle.net/11321/690>, CLARIN-PL digital repository
9. Kalchbrenner, N., Grefenstette, E., Blunsom, P.: A convolutional neural network for modelling sentences. In: Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics. pp. 655–665 (2014)
10. Kennedy, A., Inkpen, D.: Sentiment classification of movie reviews using contextual valence shifters. Computational Intelligence **22**, 110–125 (2006)
11. Kim, Y.: Convolutional neural networks for sentence classification. In: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing. pp. 1746–1751 (2014)
12. Kocoń, J., Milkowski, P., Zaśko-Zielińska, M.: Multi-level sentiment analysis of PolEmo 2.0: Extended corpus of multi-domain consumer reviews. In: Proceedings of the 23rd Conference on Computational Natural Language Learning. pp. 980–991 (2019)
13. Kumar, A., Irsoy, O., Ondruska, P., Iyyer, M., Bradbury, J., Gulrajani, I., Zhong, V., Paulus, R., Socher, R.: Ask me anything: Dynamic memory networks for natural language processing. In: Proceedings of the 33rd International Conference on International Conference on Machine Learning. vol. 48, p. 1378–1387 (2016)
14. Letarte, G., Paradis, F., Giguère, P., Laviolette, F.: Importance of self-attention for sentiment analysis. In: Proceedings of the 2018 EMNLP Workshop BlackboxNLP. pp. 267–275 (2018)
15. Lin, Z., Feng, M., dos Santos, C.N., Yu, M., Xiang, B., Zhou, B., Bengio, Y.: A structured self-attentive sentence embedding. In: 5th International Conference on Learning Representations (2017)
16. Liu, B.: Sentiment Analysis and Opinion Mining. Synthesis Lectures on Human Language Technologies, Morgan & Claypool Publishers (2012)

17. May, P.: German ELMo Model (2019), <https://github.com/t-systems-on-site-services-gmbh/german-elmo-model>, accessed: 2020-01-20
18. McCann, B., Bradbury, J., Xiong, C., Socher, R.: Learned in translation: Contextualized word vectors. In: *Advances in Neural Information Processing Systems* 30, pp. 6294–6305 (2017)
19. Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J.: Distributed representations of words and phrases and their compositionality. In: *Advances in Neural Information Processing Systems* 26, pp. 3111–3119 (2013)
20. Mohammad, S.M.: Sentiment analysis: Detecting valence, emotions, and other affectual states from text. In: *Emotion measurement*, pp. 201–237. Elsevier (2016)
21. Pang, B., Lee, L., Vaithyanathan, S.: Thumbs up? sentiment classification using machine learning techniques. In: *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing*. pp. 79–86 (2002)
22. Paulus, R., Socher, R., Manning, C.D.: Global belief recursive neural networks. In: *Advances in Neural Information Processing Systems* 27, pp. 2888–2896 (2014)
23. Pennington, J., Socher, R., Manning, C.: Glove: Global vectors for word representation. In: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*. pp. 1532–1543 (2014)
24. Peters, M., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., Zettlemoyer, L.: Deep contextualized word representations. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. pp. 2227–2237 (2018)
25. Radford, A., Narasimhan, K., Salimans, T., Sutskever, I.: Improving language understanding by generative pre-training (2018)
26. dos Santos, C.N., Tan, M., Xiang, B., Zhou, B.: Attentive pooling networks. *arXiv* (2016)
27. Shaw, P., Uszkoreit, J., Vaswani, A.: Self-attention with relative position representations. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. pp. 464–468 (2018)
28. Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C.D., Ng, A., Potts, C.: Recursive deep models for semantic compositionality over a sentiment treebank. In: *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*. pp. 1631–1642 (2013)
29. Tai, K.S., Socher, R., Manning, C.D.: Improved semantic representations from tree-structured long short-term memory networks. In: *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing*. pp. 1556–1566 (2015)
30. Turney, P.: Thumbs up or thumbs down? semantic orientation applied to unsupervised classification of reviews. In: *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*. pp. 417–424 (2002)
31. Turney, P.D., Pantel, P.: From frequency to meaning: Vector space models of semantics. *J. Artif. Intell. Res.* **37**, 141–188 (2010)
32. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: *Advances in Neural Information Processing Systems* 30: Annual Conference on Neural Information Processing Systems. pp. 5998–6008 (2017)
33. Wojatzki, M., Ruppert, E., Holschneider, S., Zesch, T., Biemann, C.: GermEval 2017: Shared Task on Aspect-based Sentiment in Social Media Customer Feedback. In: *Proceedings of the GermEval 2017*. pp. 1–12 (2017)