

Attention-based Conditioning Methods for External Knowledge Integration

Katerina Margatina¹, Christos Baziotis²*, Alexandros Potamianos^{1,3,4}

¹School of ECE, National Technical University of Athens, Athens, Greece

²School of Informatics, University of Edinburgh, UK

³Signal Analysis and Interpretation Laboratory (SAIL), USC, Los Angeles, USA

⁴Behavioral Signal Technologies, Los Angeles, USA

el12108@central.ntua.gr, c.baziotis@sms.ed.ac.uk,
potam@central.ntua.gr

Abstract

In this paper, we present a novel approach for incorporating external knowledge in Recurrent Neural Networks (RNNs). We propose the integration of lexicon features into the self-attention mechanism of RNN-based architectures. This form of conditioning on the attention distribution, enforces the contribution of the most salient words for the task at hand. We introduce three methods, namely attentional concatenation, feature-based gating and affine transformation. Experiments on six benchmark datasets show the effectiveness of our methods. Attentional feature-based gating yields consistent performance improvement across tasks. Our approach is implemented as a simple add-on module for RNN-based models with minimal computational overhead and can be adapted to any deep neural architecture.

1 Introduction

Modern deep learning algorithms often do away with feature engineering and learn latent representations directly from raw data that are given as input to Deep Neural Networks (DNNs) (Mikolov et al., 2013; McCann et al., 2017; Peters et al., 2018). However, it has been shown that linguistic knowledge (manually or semi-automatically encoded into lexicons and knowledge bases) can significantly improve DNN performance for Natural Language Processing (NLP) tasks, such as natural language inference (Mrkšić et al., 2017), language modelling (Ahn et al., 2016), named entity recognition (Ghaddar and Langlais, 2018) and relation extraction (Vashishth et al., 2018).

For NLP tasks, external sources of information are typically incorporated into deep neural architectures by processing the raw input *in the context* of such external linguistic knowledge. In machine

learning, this contextual processing is known as *conditioning*; the computation carried out by a model is conditioned or modulated by information extracted from an auxiliary input. The most commonly-used method of conditioning is concatenating a representation of the external information to the input or hidden network layers.

Attention mechanisms (Bahdanau et al., 2015; Vaswani et al., 2017; Lin et al., 2017) are a key ingredient for achieving state-of-the-art performance in tasks such as textual entailment (Rocktäschel et al., 2016), question answering (Xiong et al., 2017), and neural machine translation (Wu et al., 2016). Often task-specific attentional architectures are proposed in the literature to further improve DNN performance (Dhingra et al., 2017; Xu et al., 2015; Barrett et al., 2018).

In this work, we propose a novel way of utilizing word-level prior information encoded in linguistic, sentiment, and emotion lexicons, to improve classification performance. Usually, lexicon features are concatenated to word-level representations (Wang et al., 2016; Yang et al., 2017; Trotzek et al., 2018), as additional features to the embedding of each word or the hidden states of the model. By contrast, we propose to incorporate them into the self-attention mechanism of RNNs. Our goal is to enable the self-attention mechanism to identify the most informative words, by directly conditioning on their additional lexicon features.

Our contributions are the following: (1) we propose an alternative way for incorporating external knowledge to RNN-based architectures, (2) we present empirical results that our proposed approach consistently outperforms strong baselines, and (3) we report state-of-the-art performance in two datasets. We make our source code publicly available¹.

*The research was conducted when the author was a researcher at School of ECE, NTUA in Athens, Greece.

¹<https://github.com/mourga/affective-attention>

2 Related Work

In the traditional machine learning literature where statistical models are based on sparse features, affective lexicons have been shown to be highly effective for tasks such as sentiment analysis, as they provide additional information not captured in the raw training data (Hu and Liu, 2004; Kim and Hovy, 2004; Ding et al., 2008; Yu and Dredze, 2014; Taboada et al., 2011). After the emergence of pretrained word representations (Mikolov et al., 2013; Pennington et al., 2014), the use of lexicons is no longer common practice, since word embeddings can also capture some of the affective meaning of these words.

Recently, there have been notable contributions towards integrating linguistic knowledge into DNNs for various NLP tasks. For sentiment analysis, Teng et al. (2016) integrate lexicon features to an RNN-based model with a custom weighted-sum calculation of word features. Shin et al. (2017) propose three convolutional neural network specific methods of lexicon integration achieving state-of-the-art performance on two datasets. Kumar et al. (2018) concatenate features from a knowledge base to word representations in an attentive bidirectional LSTM architecture, also reporting state-of-the-art results. For sarcasm detection, Yang et al. (2017) incorporate psycholinguistic, stylistic, structural, and readability features by concatenating them to paragraph and document-level representations.

Furthermore, there is limited literature regarding the development and evaluation of methods for combining representations in deep neural networks. Peters et al. (2017) claim that concatenation, non-linear mapping and attention-like mechanisms are unexplored methods for including language model representations in their sequence model. They employ simple concatenation, leaving the exploration of other methods to future work. Dumoulin et al. (2018) provide an overview of feature-wise transformations such as concatenation-based conditioning, conditional biasing and gating mechanisms. They review the effectiveness of conditioning methods in tasks such as visual question answering (Strub et al., 2018), style transfer (Dumoulin et al., 2017) and language modeling (Dauphin et al., 2017). They also extend the work by Perez et al. (2017), which proposes the Feature-wise Linear Modulation (FiLM) framework, and investigate its applications in vi-

sual reasoning tasks. Balazs and Matsuo (2019) provide an empirical study showing the effects of different ways of combining character and word representations in word-level and sentence-level evaluation tasks. Some of the reported findings are that gating conditioning performs consistently better across a variety of word similarity and relatedness tasks.

3 Proposed Model

3.1 Network Architecture

Word Embedding Layer. The input sequence of words w_1, w_2, \dots, w_T is projected to a low-dimensional vector space R^W , where W is the size of the embedding layer and T the number of words in a sentence. We initialize the weights of the embedding layer with pretrained word embeddings.

LSTM Layer. A Long Short-Term Memory unit (LSTM) (Hochreiter and Schmidhuber, 1997) takes as input the words of a sentence and produces the word annotations h_1, h_2, \dots, h_T , where h_i is the hidden state of the LSTM at time-step i , summarizing all sentence information up to w_i .

Self-Attention Layer. We use a self-attention mechanism (Cheng et al., 2016) to find the relative importance of each word for the task at hand. The attention mechanism assigns a score a_i to each word annotation h_i . We compute the fixed representation r of the input sequence, as the weighted sum of all the word annotations. Formally:

$$a_i = \text{softmax}(v_a^T f(h_i)) \quad (1)$$

$$r = \sum_{i=1}^T a_i h_i \quad (2)$$

where $f(\cdot)$ corresponds to a non-linear transformation $\tanh(W_a h_i + b_a)$ and W_a, b_a, v_a are the parameters of the attention layer.

Lexicons	Annotations	# dim.	# words
LIWC	psycho-linguistic	73	18,504
Bing Liu	valence	1	2,477
AFINN	sentiment	1	6,786
MPQA	sentiment	4	6,886
SemEval15	sentiment	1	1,515
Emolex	emotion	19	14,182

Table 1: The lexicons used as external knowledge.

3.2 External Knowledge

In this work, we augment our models with existing linguistic and affective knowledge from human experts. Specifically, we leverage lexica containing psycho-linguistic, sentiment and emotion annotations. We construct a feature vector $c(w_i)$ for every word in the vocabulary by concatenating the word’s annotations from the lexicons shown in Table 1. For missing words we append zero in the corresponding dimension(s) of $c(w_i)$.

3.3 Conditional Attention Mechanism

We extend the standard self-attention mechanism (Eq. 1, 2), in order to condition the attention distribution of a given sentence, on each word’s prior lexical information. To this end, we use as input to the self-attention layer both the word annotation h_i , as well as the lexicon feature $c(w_i)$ of each word. Therefore, we replace $f(h_i)$ in Eq. 1 with $f(h_i, c(w_i))$. Specifically, we explore three conditioning methods, which are illustrated in Figure 1. We refer to the conditioning function as $f_i(\cdot)$, the weight matrix as W_i and the biases as b_i , where i is an indicative letter for each method. We present our results in Section 5 (Table 3) and we denote the three conditioning methods as “*conc.*”, “*gate*” and “*affine*” respectively.

Attentional Concatenation. In this approach, as illustrated in Fig. 1(a), we learn a function of the concatenation of each word annotation h_i with its lexicon features $c(w_i)$. The intuition is that by adding extra dimensions to h_i , learned representations are more discriminative. Concretely:

$$f_c(h_i, c(w_i)) = \tanh(W_c[h_i \parallel c(w_i)] + b_c) \quad (3)$$

where \parallel denotes the concatenation operation and W_c, b_c are learnable parameters.

Attentional Feature-based Gating. The second approach, illustrated in Fig. 1(b), learns a feature mask, which is applied on each word annotation h_i . Specifically, a gate mechanism with a sigmoid activation function, generates a mask-vector from each $c(w_i)$ with values between 0 and 1 (black and white dots in Fig. 1(b)). Intuitively, this gating mechanism selects salient dimensions (i.e. features) of h_i , conditioned on the lexical information. Formally:

$$f_g(h_i, c(w_i)) = \sigma(W_g c(w_i) + b_g) \odot h_i \quad (4)$$

where \odot denotes element-wise multiplication and W_g, b_g are learnable parameters.

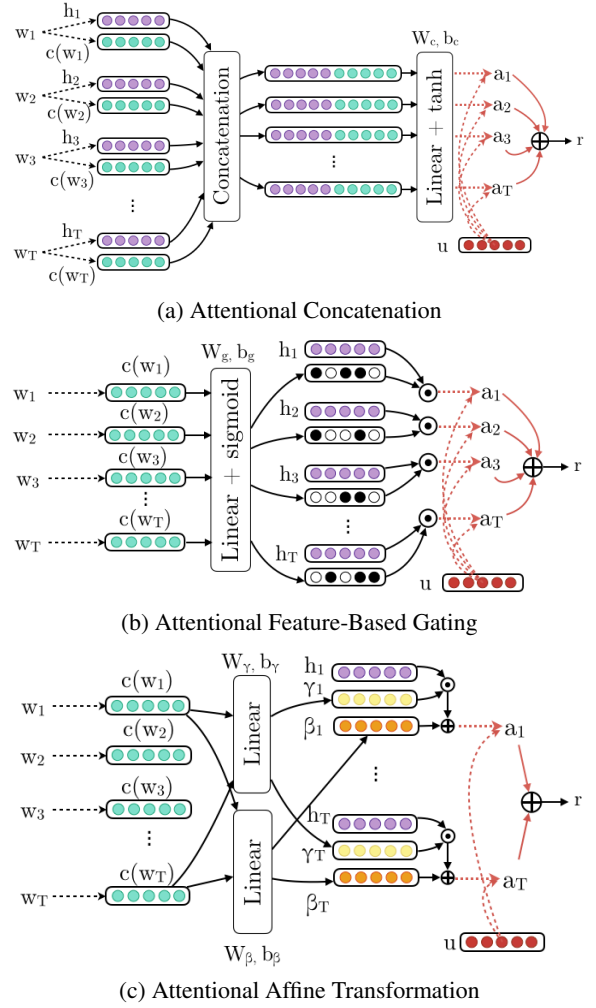


Figure 1: The proposed conditioning methods of the self-attention mechanism.

Attentional Affine Transformation. The third approach, shown in Fig. 1(c), is adopted from the work of Perez et al. (2017) and applies a feature-wise affine transformation to the latent space of the hidden states. Specifically, we use the lexicon features $c(w_i)$, in order to conditionally generate the corresponding scaling $\gamma(\cdot)$ and shifting $\beta(\cdot)$ vectors. Concretely:

$$f_a(h_i, c(w_i)) = \gamma(c(w_i)) \odot h_i + \beta(c(w_i)) \quad (5)$$

$$\gamma(x) = W_\gamma x + b_\gamma, \quad \beta(x) = W_\beta x + b_\beta \quad (6)$$

where $W_\gamma, W_\beta, b_\gamma, b_\beta$ are learnable parameters.

3.4 Baselines

We employ two baselines: The first baseline is an LSTM-based architecture augmented with a self-attention mechanism (Sec. 3.1) with no external knowledge. The second baseline incorporates lexicon information by concatenating the $c(w_i)$ vec-

Dataset	Study	Task	Domain	Classes	N_{train}	N_{test}
SST-5	Socher et al. (2013)	Sentiment	Movie Reviews	5	9,645	2,210
Sent17	Rosenthal et al. (2017)	Sentiment	Twitter	3	49,570	12,284
PhychExp	Wallbott and Scherer (1986)	Emotion	Experiences	7	1000	6480
Irony18	Van Hee et al. (2018)	Irony	Twitter	4	3,834	784
SCv1	Lukin and Walker (2013)	Sarcasm	Debate Forums	2	1000	995
SCv2	Oraby et al. (2016)	Sarcasm	Debate Forums	2	1000	2260

Table 2: Description of benchmark datasets. We split 10% of the train set to serve as the validation set.

Model	SST-5	Sent17	PhychExp	Irony18	SCv1	SCv2
baseline	43.5 ± 0.5	68.3 ± 0.2	53.2 ± 0.8	46.3 ± 1.4	64.1 ± 0.5	74.0 ± 0.7
emb. conc.	43.3 ± 0.6	68.4 ± 0.2	57.1 ± 1.2	48.1 ± 1.2	64.2 ± 0.7	74.2 ± 0.7
conc.	44.0 ± 0.7	68.6 ± 0.3	54.3 ± 0.6	47.4 ± 0.9	65.1 ± 0.6	74.3 ± 1.2
gate	44.2 ± 0.4	68.7 ± 0.3	53.4 ± 1.0	48.5 ± 0.7	64.7 ± 0.7	74.3 ± 1.2
affine	43.2 ± 0.7	68.5 ± 0.3	53.1 ± 0.9	45.3 ± 1.5	60.3 ± 0.8	74.0 ± 1.0
gate+emb.conc.	46.2 ± 0.5	68.9 ± 0.3	57.2 ± 1.1	48.4 ± 1.0	64.9 ± 0.6	74.4 ± 0.9
state-of-the-art	51.7	68.5	57.0	53.6	69.0	76.0
	Shen et al. (2018) Cliche (2017) Felbo et al. (2017) Baziotis et al. (2018) Felbo et al. (2017) Ilić et al. (2018)					

Table 3: Comparison across benchmark datasets. Reported values are averaged across ten runs. All reported measures are F_1 scores, apart from *SST* – 5 which is evaluated with *Accuracy*.

tors to the word representations in the embedding layer. In Table 3 we use the abbreviations “*baseline*” and “*emb. conc.*” for the two baseline models respectively.

4 Experiments

Lexicon Features. As prior knowledge, we leverage the lexicons presented in Table 1. We selected widely-used lexicons that represent different facets of affective and psycho-linguistic features, namely; LIWC (Tausczik and Pennebaker, 2010), Bing Liu Opinion Lexicon (Hu and Liu, 2004), AFINN (Nielsen, 2011), Subjectivity Lexicon (Wilson et al., 2005), SemEval 2015 English Twitter Lexicon (Svetlana Kiritchenko and Mohammad, 2014), and NRC Emotion Lexicon (EmoLex) (Mohammad and Turney, 2013).

Datasets. The proposed framework can be applied to different domains and tasks. In this paper, we experiment with sentiment analysis, emotion recognition, irony, and sarcasm detection. Details of the benchmark datasets are shown in Table 2.

Preprocessing. To preprocess the words, we use the tool *Ekphrasis* (Baziotis et al., 2017). After tokenization, we map each word to the corresponding pretrained word representation: Twitter-specific word2vec embeddings (Chronopoulou

et al., 2018) for the Twitter datasets, and fast-text (Bojanowski et al., 2017) for the rest.

Experimental Setup. For all methods, we employ a single-layer LSTM model with 300 neurons augmented with a self-attention mechanism, as described in Section 3. As regularization techniques, we apply early stopping, Gaussian noise $N(0, 0.1)$ to the word embedding layer, and dropout to the LSTM layer with $p = 0.2$. We use Adam to optimize our networks (Kingma and Ba, 2014) with mini-batches of size 64 and clip the norm of the gradients (Pascanu et al., 2013) at 0.5, as an extra safety measure against exploding gradients. We also use PyTorch (Paszke et al., 2017) and scikit-learn (Pedregosa et al., 2011).

5 Results & Analysis

We compare the performance of the three proposed conditioning methods with the two baselines and the state-of-the-art in Table 3. We also provide results for the combination of our best method, attentional feature-based gating, and the second baseline model (*emb. conc.*).

The results show that incorporating external knowledge in RNN-based architectures consistently improves performance over the baseline for all datasets. Furthermore, feature-based gating im-

when	a	mistake	occurred	at	work	which	i	was	not	responsible	for	.	this	was	disclosed	later	.	anger	✗
0.055	0.048	0.053	0.058	0.049	0.048	0.046	0.043	0.044	0.048	0.106	0.071	0.065	0.059	0.052	0.054	0.050	0.052		
when	a	mistake	occurred	at	work	which	i	was	not	responsible	for	.	this	was	disclosed	later	.	guilt	✓
0.022	0.023	0.153	0.150	0.138	0.133	0.068	0.025	0.022	0.033	0.038	0.029	0.035	0.027	0.023	0.026	0.024	0.029		

Figure 2: Attention heatmap of a *PsychExp* random test sample. The first attention distribution is created with the *baseline* model without lexicon feature integration, while the second with the combination of our attentional feature-based gating method and the concatenation to word embeddings baseline (*gate+emb.conc.*).

proves upon baseline concatenation in the embedding layer across benchmarks, with the exception of *PsychExp* dataset.

For the *Sent17* dataset we achieve state-of-the-art F_1 score using the feature-based gating method; we further improve performance when combining gating with the *emb. conc.* method. For *SST-5*, we observe a significant performance boost with combined attentional gating and embedding conditioning (*gate + emb. conc.*). For *PsychExp*, we marginally outperform the state-of-the-art also with the combined method, while for *Irony18*, feature-based gating yields the best results. Finally, concatenation based conditioning is the top method for *SCv1*, and the combination method for *SCv2*.

Overall, attentional feature-based gating is the best performing conditioning method followed by concatenation. Attentional affine transformation underperforms, especially, for smaller datasets; this is probably due to the higher capacity of this model. This is particularly interesting since gating (Eq. 4) is a special case of the affine transformation method (Eq. 5), where the shifting vector β is zero and the scaling vector γ is bounded to the range $[0, 1]$ (Eq. 6). Interestingly, the combination of gating with traditional embedding-layer concatenation gives additional performance gains for most tasks, indicating that there are synergies to exploit in various conditioning methods.

In addition to the performance improvements, we can visually evaluate the effect of conditioning the attention distribution on prior knowledge and improve the interpretability of our approach. As we can see in Figure 2, lexicon features guide the model to attend to more salient words and thus predict the correct class.

6 Conclusions & Future work


We introduce three novel attention-based conditioning methods and compare their effectiveness

with traditional concatenation-based conditioning. Our methods are simple, yet effective, achieving consistent performance improvement for all datasets. Our approach can be applied to any RNN-based architecture as a extra module to further improve performance with minimal computational overhead.

In the future, we aim to incorporate more elaborate linguistic resources (e.g. knowledge bases) and to investigate the performance of our methods on more complex NLP tasks, such as named entity recognition and sequence labelling, where prior knowledge integration is an active area of research.

Acknowledgements

We would like to thank our colleagues Alexandra Chronopoulou and Georgios Paraskevopoulos for their helpful suggestions and comments. This work has been partially supported by computational timegranted from the Greek Research & Technology Network (GR-NET) in the National HPC facility - ARIS. We thank NVIDIA for supporting this work by donating a TitanX GPU.

 This work was conducted within the scope of the Research and Innovation Action *Gourmet*, which has received funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement No 825299.

References

- Sungjin Ahn, Heeyoul Choi, Tanel Pärnamaa, and Yoshua Bengio. 2016. [A neural knowledge language model](#). *arXiv preprint arXiv:1608.00318*.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. [Neural machine translation by jointly learning to align and translate](#). In *Proceedings of the International Conference on Learning Representations*.

- Jorge Balazs and Yutaka Matsuo. 2019. [Gating mechanisms for combining character and word-level word representations: an empirical study](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop*, pages 110–124, Minneapolis, Minnesota. Association for Computational Linguistics.
- Maria Barrett, Joachim Bingel, Nora Hollenstein, Marek Rei, and Anders Søgaard. 2018. [Sequence classification with human attention](#). In *Proceedings of the Conference on Computational Natural Language Learning*, pages 302–312.
- Christos Baziotis, Athanasia Nikolaos, Pinelopi Papalampidi, Athanasia Kolovou, Georgios Paraskevopoulos, Nikolaos Ellinas, and Alexandros Potamianos. 2018. [Ntua-slp at semeval-2018 task 3: Tracking ironic tweets using ensembles of word and character level attentive rnns](#). In *Proceedings of the International Workshop on Semantic Evaluation*, pages 613–621.
- Christos Baziotis, Nikos Pelekis, and Christos Douk-eridis. 2017. [Datastories at semeval-2017 task 4: Deep lstm with attention for message-level and topic-based sentiment analysis](#). In *Proceedings of the International Workshop on Semantic Evaluation*, pages 747–754.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. [Enriching word vectors with subword information](#). *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Jianpeng Cheng, Li Dong, and Mirella Lapata. 2016. [Long short-term memory-networks for machine reading](#). In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 551–561.
- Alexandra Chronopoulou, Aikaterini Margatina, Christos Baziotis, and Alexandros Potamianos. 2018. [Ntua-slp at irst 2018: Ensemble of neural transfer methods for implicit emotion classification](#). In *Proceedings of the Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 57–64.
- Mathieu Cliche. 2017. [Bb_twtr at semeval-2017 task 4: Twitter sentiment analysis with cnns and lstms](#). In *Proceedings of the International Workshop on Semantic Evaluation*, pages 573–580.
- Yann N. Dauphin, Angela Fan, Michael Auli, and David Grangier. 2017. [Language modeling with gated convolutional networks](#). In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML'17*, pages 933–941. JMLR.org.
- Bhuvan Dhingra, Hanxiao Liu, Zhilin Yang, William Cohen, and Ruslan Salakhutdinov. 2017. [Gated-attention readers for text comprehension](#). In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, pages 1832–1846.
- Xiaowen Ding, Bing Liu, and Philip S. Yu. 2008. [A holistic lexicon-based approach to opinion mining](#). In *Proceedings of the International Conference on Web Search and Data Mining*, pages 231–240.
- Vincent Dumoulin, Ethan Perez, Nathan Schucher, Florian Strub, Harm de Vries, Aaron Courville, and Yoshua Bengio. 2018. [Feature-wise transformations](#). *Distill*.
- Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur. 2017. [A learned representation for artistic style](#). *International Conference on Learning Representations*.
- Bjarke Felbo, Alan Mislove, Anders Søgaard, Iyad Rahwan, and Sune Lehmann. 2017. [Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm](#). In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 1615–1625.
- Abbas Ghaddar and Phillippe Langlais. 2018. [Robust lexical features for improved neural network named-entity recognition](#). In *Proceedings of the International Conference on Computational Linguistics*, pages 1896–1907.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. [Long short-term memory](#). *Neural Computation*, 9:1735–1780.
- Minqing Hu and Bing Liu. 2004. [Mining and summarizing customer reviews](#). In *Proceedings of the International Conference on Knowledge Discovery and Data Mining*, pages 168–177.
- Suzana Ilić, Edison Marrese-Taylor, Jorge Balazs, and Yutaka Matsuo. 2018. [Deep contextualized word representations for detecting sarcasm and irony](#). In *Proceedings of the Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 2–7.
- Soo-Min Kim and Eduard H. Hovy. 2004. [Determining the sentiment of opinions](#). In *Proceedings of the International Conference on Computational Linguistics*.
- Diederik Kingma and Jimmy Ba. 2014. [Adam: A method for stochastic optimization](#). *arXiv preprint arXiv:1412.6980*.
- Abhishek Kumar, Daisuke Kawahara, and Sadao Kurohashi. 2018. [Knowledge-enriched two-layered attention network for sentiment analysis](#). In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics*, pages 253–258.
- Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. 2017. [A structured self-attentive sentence embedding](#). In *Proceedings of the International Conference on Learning Representations*.

- Stephanie Lukin and Marilyn Walker. 2013. [Really? well. apparently bootstrapping improves the performance of sarcasm and nastiness classifiers for online dialogue.](#) In *Proceedings of the Workshop on Language Analysis in Social Media*, pages 30–40.
- Bryan McCann, James Bradbury, Caiming Xiong, and Richard Socher. 2017. [Learned in translation: Contextualized word vectors.](#) In *Advances in Neural Information Processing Systems*, pages 6294–6305.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. [Distributed representations of words and phrases and their compositionality.](#) In *Proceedings of the International Conference on Neural Information Processing Systems*, pages 3111–3119.
- Saif M. Mohammad and Peter D. Turney. 2013. [Crowdsourcing a word-emotion association lexicon.](#) 29(3):436–465.
- Nikola Mrkšić, Ivan Vulić, Diarmuid Ó Séaghdha, Ira Leviant, Roi Reichart, Milica Gašić, Anna Korhonen, and Steve Young. 2017. [Semantic specialization of distributional word vector spaces using monolingual and cross-lingual constraints.](#) 5:309–324.
- F. Å. Nielsen. 2011. [Afinn.](#)
- Shereen Oraby, Vrindavan Harrison, Lena Reed, Ernesto Hernandez, Ellen Riloff, and Marilyn Walker. 2016. [Creating and characterizing a diverse corpus of sarcasm in dialogue.](#) In *Proceedings of the Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 31–41.
- Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. 2013. [On the difficulty of training recurrent neural networks.](#) *International Conference on Machine Learning*, 28:1310–1318.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. [Automatic differentiation in pytorch.](#)
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, and others. 2011. [Scikit-learn: Machine learning in Python.](#) *Journal of Machine Learning Research*, 12:2825–2830.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. [Glove: Global vectors for word representation.](#) In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 1532–1543.
- Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron C. Courville. 2017. [Film: Visual reasoning with a general conditioning layer.](#) *CoRR*, abs/1709.07871.
- Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. [Deep contextualized word representations.](#) In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics*, pages 2227–2237.
- Matthew E. Peters, Waleed Ammar, Chandra Bhagavatula, and Russell Power. 2017. [Semi-supervised sequence tagging with bidirectional language models.](#) *CoRR*, abs/1705.00108.
- Tim Rocktäschel, Edward Grefenstette, Karl Moritz Hermann, Tomáš Kociský, and Phil Blunsom. 2016. [Reasoning about entailment with neural attention.](#) In *Proceedings of the International Conference on Learning Representations*.
- Sara Rosenthal, Noura Farra, and Preslav Nakov. 2017. [Semeval-2017 task 4: Sentiment analysis in twitter.](#) In *Proceedings of the International Workshop on Semantic Evaluation*, pages 502–518.
- Tao Shen, Tianyi Zhou, Guodong Long, Jing Jiang, Shirui Pan, and Chengqi Zhang. 2018. [Disan: Directional self-attention network for rnn/cnn-free language understanding.](#) In *Association for the Advancement of Artificial Intelligence*.
- Bonggun Shin, Timothy Lee, and Jinho D. Choi. 2017. [Lexicon integrated cnn models with attention for sentiment analysis.](#) In *Proceedings of the Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 149–158.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. [Recursive deep models for semantic compositionality over a sentiment treebank.](#) In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642.
- Florian Strub, Mathieu Seurin, Ethan Perez, Harm de Vries, Jeremie Mary, Philippe Preux, and Aaron Courville/Olivier Pietquin. 2018. [Visual reasoning with multi-hop feature modulation.](#) In *The European Conference on Computer Vision (ECCV)*.
- Xiaodan Zhu Svetlana Kiritchenko and Saif M. Mohammad. 2014. [Sentiment analysis of short informal texts.](#) 50:723–762.
- Maite Taboada, Julian Brooke, Milan Tofiloski, Kimberly Voll, and Manfred Stede. 2011. [Lexicon-based methods for sentiment analysis.](#) *Comput. Linguist.*, 37(2):267–307.
- Yla R. Tausczik and James W. Pennebaker. 2010. [The psychological meaning of words: Liwc and computerized text analysis methods.](#) *Journal of Language and Social Psychology*, 29(1):24–54.
- Zhiyang Teng, Duy-Tin Vo, and Yue Zhang. 2016. [Context-sensitive lexicon features for neural sentiment analysis.](#) In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language*

- Processing*, pages 1629–1638, Austin, Texas. Association for Computational Linguistics.
- Marcel Trotzek, Sven Koitka, and Christoph M. Friedrich. 2018. [Utilizing neural networks and linguistic metadata for early detection of depression indications in text sequences](#). *CoRR*, abs/1804.07000.
- Cynthia Van Hee, Els Lefever, and Veronique Hoste. 2018. [Semeval-2018 task 3: Irony detection in english tweets](#). In *Proceedings of The International Workshop on Semantic Evaluation*, pages 39–50.
- Shikhar Vashishth, Rishabh Joshi, Sai Suman Prayaga, Chiranjib Bhattacharyya, and Partha Talukdar. 2018. [Reside: Improving distantly-supervised neural relation extraction using side information](#). In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 1257–1266.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems*, pages 5998–6008.
- Harald G. Wallbott and Klaus R. Scherer. 1986. [How universal and specific is emotional experience? evidence from 27 countries on five continents](#). *Information (International Social Science Council)*, 25:763–795.
- Bingning Wang, Shangmin Guo, Kang Liu, Shizhu He, and Jun Zhao. 2016. [Employing external rich knowledge for machine comprehension](#). In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI'16*, pages 2929–2935. AAAI Press.
- Theresa Wilson, Paul Hoffmann, Swapna Somasundaran, Jason Kessler, Janyce Wiebe, Yejin Choi, Claire Cardie, Ellen Riloff, and Siddharth Patwardhan. 2005. [Opinionfinder: A system for subjectivity analysis](#). In *Proceedings of HLT/EMNLP on Interactive Demonstrations*. Association for Computational Linguistics.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2016. [Google’s neural machine translation system: Bridging the gap between human and machine translation](#). *CoRR*, abs/1609.08144.
- Caiming Xiong, Victor Zhong, and Richard Socher. 2017. [Dynamic coattention networks for question answering](#). In *International Conference on Learning Representations*.
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. [Show, attend and tell: Neural image caption generation with visual attention](#). In *Proceedings of the International Conference on Machine Learning*, pages 2048–2057.
- Fan Yang, Arjun Mukherjee, and Eduard Dragut. 2017. [Satirical news detection and analysis using attention mechanism and linguistic features](#). In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 1979–1989.
- Mo Yu and Mark Dredze. 2014. [Improving lexical embeddings with semantic knowledge](#). In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, pages 545–550.