

Attention Word Embedding

Shashank Sonkar
Rice University
ss164@rice.edu

Andrew E. Waters
Rice University
aew2@rice.edu

Richard G. Baraniuk
Rice University
richb@rice.edu

Abstract

Word embedding models learn semantically rich vector representations of words and are widely used to initialize natural processing language (NLP) models. The popular continuous bag-of-words (CBOW) model of word2vec learns a vector embedding by masking a given word in a sentence and then using the other words as a context to predict it. A limitation of CBOW is that it equally weights the context words when making a prediction, which is inefficient, since some words have higher predictive value than others. We tackle this inefficiency by introducing the *Attention Word Embedding* (AWE) model, which integrates the attention mechanism into the CBOW model. We also propose AWE-S, which incorporates subword information. We demonstrate that AWE and AWE-S outperform the state-of-the-art word embedding models both on a variety of word similarity datasets and when used for initialization of NLP models.

1 Introduction

Word embedding models learn vector representations of words such that words that are semantically related are close to each other in the vector space. Popular word embedding models include word2vec (Mikolov et al., 2013a; Mikolov et al., 2013b), GloVe (Pennington et al., 2014), and fastText (Bojanowski et al., 2017). Word embedding models are often used to initialize deep learning models for various natural language processing (NLP) tasks such as machine translation (Bahdanau et al., 2015), part-of-speech tagging (Ling et al., 2015a), and sentiment analysis (Wang et al., 2016), leading to improved performance over training the models from scratch.

The core idea underlying the training methodology of all the above embedding models is that “a word is characterized by the company it keeps” (Firth, 1957). For instance, the continuous bag-of-words (CBOW) model of word2vec predicts a randomly selected word in a sentence using the other words in the sentence as a context. The context words are typically treated equally by these models. However, clearly some words in the context will be more predictive of the masked word than others, and intuitively, should be weighted more heavily.

A promising remedy that treats each context word differently is the *attention mechanism* (Vaswani et al., 2017), which enables learning the relative importance of input features for the prediction of the desired output. Attention has become the de facto standard and state-of-the-art in modern NLP for a range of different tasks (Zhang et al., 2019; Devlin et al., 2019; Dai et al., 2019). However, to date, there have been only a few attempts to employ it to improve the performance and interpretability of word embedding models (Ling et al., 2015b; Schick and Schütze, 2019).

1.1 Contributions

In this paper, we propose the *Attention Word Embedding* (AWE), a new word embedding model that integrates the attention mechanism into the CBOW model of word2vec.

AWE consists of two components. First, AWE uses a variant of self-attention (Vaswani et al., 2017) to narrow down relevant words in the context for the prediction of the masked word. In contrast to prior

work (Ling et al., 2015b; Schick and Schütze, 2019), the attention weights are a function of both the context words and the target/masked word. Second, AWE embeds the context words and the masked word representations in a shared subspace, since both representations embody the meaning of the word. Note that this is in sharp contrast with CBOW, which learns two different embeddings for each word, one when the word is present in the context and another when it is masked.

We also introduce a variant of AWE, AWE-S, which leverages subword information to enrich the word vectors. The idea to use subwords is inspired from fastText (Bojanowski et al., 2017), which defines subwords of a word as its character n-grams. Incorporating subword information is particularly valuable for improving word vector representations for rare and unseen words.

Our experimental evaluations show the superior performance of AWE and AWE-S than many existing word embedding models in both language modeling tasks and a number of downstream NLP applications including natural language inference, sentence semantic relatedness, and paraphrase detection. Furthermore, we analyze and interpret the role of the attention weights in AWE, which provide interesting insights into the workings of the attention mechanism.

2 Preliminaries

2.1 Word2Vec

Word2vec is one of the standard word embedding models (Mikolov et al., 2013a; Mikolov et al., 2013b). There are two architectures proposed for word2vec: CBOW and Skip-Gram. We focus on the CBOW model in this work.

CBOW predicts a masked word using its context (*fill in the blanks* model). For each word, it learns two vectors to represent the two roles that a word can perform - first, when it is present in the context of the masked word, and second, when it is masked.

Let $U = [\mathbf{u}_1, \dots, \mathbf{u}_N]^T \in \mathbb{R}^{N \times D}$ and $V = [\mathbf{v}_1, \dots, \mathbf{v}_N]^T \in \mathbb{R}^{N \times D}$ where N is the size of the vocabulary, and D is the size of the word vector. U models the first role and is used to calculate the context vector, \mathbf{c} , given by

$$\mathbf{c} = \sum_{i \in [-b, b] - \{0\}} \mathbf{u}_{w_i}, \quad (1)$$

where b is the size of the context window and w_i is the index of each word (w_0 is the index of the masked word; the rest are the indices of the context words). V models the second role and learns the masked word vector for w_0 , given by \mathbf{v}_{w_0} . The probability p of w_0 to occur in the context of $\{w_{-b}, \dots, w_{-1}, w_1, \dots, w_b\}$ is given by

$$p(w_0 | W_{[-b, b] - \{0\}}) \propto \exp \mathbf{v}_{w_0}^T \mathbf{c}. \quad (2)$$

The Skip-Gram model of word2vec is similar to CBOW. The key difference is that it predicts the context words using the masked word.

3 New Model: Attention Word Embedding (AWE)

3.1 The AWE Model

We now explain AWE, our proposed word embedding model that incorporates the attention mechanism. AWE augments the CBOW model of word2vec with the attention mechanism in two different ways. First, we introduce two new matrices, a key matrix, $K \in \mathbb{R}^{N \times D'}$ and a query matrix, $Q \in \mathbb{R}^{N \times D'}$, where N is the size of the vocabulary and $D' \in \mathbb{Z}$. With the attention mechanism, the context vector \mathbf{c} is not modeled as a simple sum in (1) but rather as a *weighted* sum of context word vector embeddings

$$\mathbf{c}_{\text{attn}} = \sum_{i \in [-b, b] - \{0\}} a_{w_i} \mathbf{u}_{w_i},$$

where a_{w_i} is the attention weight of each context word vector \mathbf{u}_{w_i} calculated using the key matrix K and the query matrix Q

$$a_{w_i} = \exp(\mathbf{k}_{w_0}^T \mathbf{q}_{w_i}),$$

Second, we share the weights between the context word embedding matrix and the masked word embedding matrix in our model, i.e., we set $U = V$. Sharing weights is a natural and intuitive choice, since both matrices embed the meaning of a word in its vector representation, and the meaning of the word remains the same irrespective of it occurring in the context window, or as the masked word. In CBOW, the choice is to have two separate matrices, U and V , is justified as it leads to increase in performance. However, in AWE, the intuitive choice to have U same as V works better and adds interpretability to the model. On top of that, it has an added advantage. The number of parameters in AWE are much less as compared to CBOW even though AWE has one more matrix than CBOW. The reason is that the key matrix $K \in \mathbb{R}^{N \times D'}$ and the query matrix $Q \in \mathbb{R}^{N \times D'}$ are much smaller as compared to the value matrix, $V = [\mathbf{u}_1, \dots, \mathbf{u}_D]^T \in \mathbb{R}^{N \times D}$. In our experiments, we set $D' = 50$ and $D = 500$.

3.1.1 Weight Initialization Techniques

Weight initialization is critical to the training and performance of deep neural networks (Sutskever et al., 2013). Several methods have been proposed for initializing the weight matrices which include random initialization, standard normal initialization, Xavier initialization (Glorot and Bengio, 2010), and Kaiming initialization (He et al., 2015).

We propose a new initialization scheme for initializing the key and the query matrices. The weights for both the matrices are drawn according to:

$$k_{ij}, q_{ij} \sim \mathcal{N}\left(\frac{1}{\sqrt{D'}}, 0.01\right),$$

where \mathcal{N} is the normal distribution. It is easy to verify that for any masked word, say w_m , and any word in the context of the masked word, say w_c , the initial attention weight, given by $\mathbf{k}_{w_m}^T \mathbf{q}_{w_c}$, is centered around 1. Thus, in the beginning AWE mimics the CBOW method of word2vec which assigns equal weight to each context word. As the training progresses, AWE learns the relevance of each of the context words for the prediction of the masked word. For the value matrix, we initialize the weights using $\mathcal{N}(0, 0.1)$.

3.2 AWE-S: A Subword Variant of AWE

We also propose AWE-S, a variant of AWE that leverages *subword* information. A subword of a word is a part of the word, e.g., its character n-grams or its lemmas. Subword information has been used in a number of word embedding models including fastText (Bojanowski et al., 2017) and BERT (Devlin et al., 2019) to improve the expressiveness of rare word representations. As an illustration, learning a robust representation for a word like *awing*, which occurs rarely in a corpus, is challenging due to a dearth of examples. To circumvent the scarcity of samples, a word embedding model can exploit the fact that *awe* is the verb lemma form of *awing* to learn a better embedding for *awing*.

In AWE-S, the embedding of each word is given by the sum of embeddings of all of its subwords, regardless of whether the word is a context word or the masked word in the prediction task, i.e.,

$$\tilde{\mathbf{u}}_w = \sum_{e_j \in S_w} \mathbf{u}_{e_j},$$

where S_w is the subword set of the word w . The context vector \mathbf{c} in AWE-S is given by

$$\mathbf{c}_{\text{attn}} = \sum_{i \in [-b, b] - \{0\}} a_{w_i} \left(\sum_{e_j \in S_{w_i}} \tilde{\mathbf{u}}_{e_j} \right),$$

and the probability of the masked word w_0 to occur in the context of $\{w_{-b}, \dots, w_{-1}, w_1, \dots, w_b\}$ is given by

$$p(w_0 | W_{[-b, b] - \{0\}}) \propto \exp \tilde{\mathbf{u}}_{w_0}^T \mathbf{c}_{\text{attn}}. \quad (3)$$

Note that in fastText, the context vector of (1) and the probability of w_0 to occur in the context of $\{w_{-b}, \dots, w_{-1}, w_1, \dots, w_b\}$ is given by

$$\mathbf{c} = \sum_{i \in [-b, b] - \{0\}} \sum_{e_j \in S_{w_i}} \mathbf{u}_{e_j}, \quad p(w_0 | W_{[-b, b] - \{0\}}) \propto \exp \mathbf{v}_{w_0}^T \mathbf{c}. \quad (4)$$

Model	MEN	WS353	WS353R	WS353S	SimLex999	RW(RareWords)	RG65	MTurk
AWE	0.771	0.672	0.609	0.773	0.377	0.440	0.836	0.683
CBOW	0.717	0.591	0.478	0.720	0.383	0.396	0.791	0.659
SG	0.738	0.664	0.607	0.729	0.371	0.440	0.769	0.662
GloVe	0.695	0.553	0.492	0.683	0.324	0.340	0.763	0.621

Table 1: Spearman correlation on word similarity datasets for AWE, CBOW, Skip-Gram, and GloVe.

Model	MEN	WS353	WS353R	WS353S	SimLex999	RW(RareWords)	RG65	MTurk
AWE-S	0.771	0.675	0.610	0.769	0.386	0.463	0.819	0.692
FastText	0.761	0.691	0.642	0.745	0.385	0.457	0.795	0.685

Table 2: Spearman correlation on word similarity datasets for FastText and AWE-S, which use additional subword information.

Model	SNLI (Acc)	SICK-E (Acc)	SICK-R (Pearson)	STS Benchmark (Pearson)	MRPC (F1)
AWE	65.23*	76.05*	0.783*	0.648*	81.70[†]
AWE-S	65.21*	76.05*	0.785*	0.651*	81.49 [†]
CBOW	65.07*	76.23*	0.781*	0.600*	81.26 [†]
SG	64.45*	75.04*	0.780*	0.632*	81.12 [†]
FastText	64.98*	75.06*	0.782*	0.586*	80.88 [†]
GloVe	64.18*	75.91*	0.782*	0.643*	80.60 [†]

Table 3: Evaluation results of AWE and its subword variant, AWE-S, against other word embedding methods for initializing models for downstream tasks. Accuracy is reported for SNLI dataset and SICK-E, Pearson correlation is reported for SICK-R and STS Benchmark, and F1 score is reported for MRPC. ‘*’, ‘**’ and ‘[†]’ denote accuracy, pearson correlation, and F1 score respectively.

The computations in AWE-S may appear to be similar those in fastText. However, observe that, in fastText, a word is not represented as the sum of the embeddings of its subwords when it is masked. Also, contrary to fastText, which uses character n-grams to construct the subword set for a word, AWE-S uses the noun, adjective, and verb lemmas of the word to construct the subword set. This limits the size of the subword set significantly as compared to using character n-grams, which helps to reduce the number of learnable parameters in AWE-S.

4 Experiments

We demonstrate the superior performance of AWE and AWE-S against CBOW, Skip-Gram, GloVe, and fastText on a variety of datasets, which broadly fall into two main categories:

1. **Word similarity tasks.** We use this task to evaluate the performance of word embedding models themselves. The datasets for this task contain word pairs and a semantic similarity score associated with the pairs. The scores were annotated by human subjects. These are the most widely used evaluation methods for word embedding models (Bakarov, 2018). We test AWE’s performance using spearman correlation score on eight such datasets - MEN (Bruni et al., 2014), WS353 (Finkelstein et al., 2001), WS353R (Agirre et al., 2009), WS353S (Agirre et al., 2009), SimLex999 (Hill et al., 2015), RW(RareWords) (Luong et al., 2013), RG65 (Rubenstein and Goodenough, 1965), and MTurk (Radinsky et al., 2011).
2. **Downstream NLP tasks.** These tasks evaluate performance on important NLP applications including natural language inference, semantic entailment, semantic relatedness, and paraphrase detection. Because the conventional application of word embedding models has been to initialize NLP models, we use these tasks to assess the quality of our word embeddings for initializing NLP models for downstream applications. Datasets include SNLI (Bowman et al., 2015), SICK-E (Marelli et al.,

2014), SICK-R (Marelli et al., 2014), STS Benchmark (Cer et al., 2017), and MRPC (Dolan et al., 2004). Unlike datasets for the previous task, all the datasets here have a training set.

4.1 Experimental Setup

We train AWE, CBOW, Skip-Gram, GloVe, and fastText on the wikipedia dataset (~ 17 GB in size), with vocabulary consisting of one million words. Dimensions of word embeddings for CBOW, Skip-Gram, GloVe, and fastText is 500. For AWE and AWE-S, $K, Q \in \mathbb{R}^{N \times 50}$, and $V \in \mathbb{R}^{N \times 500}$, where N is the size of the vocabulary. Widely-used parameter settings (context window size is 5, number of negative samples is 5, number of training epochs is 5) are used to train CBOW, Skip-Gram, and fastText. AWE and AWE-S are also trained for the same number of epochs. GloVe is trained for 100 epochs and parameters are set to those recommended in the paper (x-max is 100 and context window size is 10). We used open-source word embedding evaluation tool kits for assessing the quality of the models (Jastrzebski et al., 2017; Conneau and Kiela, 2018). The code for for this work is available on <https://github.com/luffycodes/attention-word-embedding.git>. The codebase is inspired from (Mai et al., 2019).

4.2 Numerical Results

We report the performance statistics of AWE, AWE-S, and the competing word embedding models on word similarity datasets and downstream NLP tasks in Tables 1, 2, and 3.

Table 1 shows performance of AWE, CBOW, Skip-Gram, and GloVe on word similarity datasets. We use spearman correlation between cosine similarity and human-annotated relatedness scores as the metric to measure the performance. AWE performs significantly better on all of the datasets except SimLex999.

Table 2 compares performance of AWE-S vs fastText. Both algorithms incorporate additional subword information into their architecture. FastText uses character n-grams, while AWE-S uses the lemma forms of the words. AWE wins against FastText on six out of the eight datasets.

Table 3 reports performance of AWE and AWE-S against CBOW, Skip-Gram, GloVe, and fastText on supervised downstream tasks. A logistic regression classifier is trained to classify sentences using pre-trained word embeddings. For these initialization applications, AWE provides improvement in performance as compared to other models over all datasets except SICK-E.

4.3 Interpretation of the attention mechanism

Studies have shown that the attention mechanism helps a neural network to attend to relevant features in the input (Wiegrefe and Pinter, 2019; Vashishth et al., 2019). Thus, motivating the integration of the attention mechanism in CBOW. However, some studies argue otherwise and show that the attention weights are poor indicators of feature importance (Jain and Wallace, 2019; Serrano and Smith, 2019). This differing opinion thus necessitates an investigation as to why attention works in the case of AWE. We visually analyze the attention weights to investigate if the attention mechanism models the importance of a context word for the masked word prediction in AWE. The investigation revealed two key findings.

First, no matter the masked word, the attention weight of the masked word and a highly frequent word is typically high. In Table 4 and 5, the attention weights for highly frequent words like *for*, *and*, *the*, *has*, and *a* (words that have cells with dark gray background) are quite high as compared to other words in the sentence. Note that even though the attention weight is large, the similarity between word vectors is very close to zero, thus not affecting the probability of prediction of the masked word.

Second, if we leave out these highly frequent words from the set of context words, we observe that the attention weights focus on more informative words in the context for the prediction of the masked word. In Table 4, for each masked word, the attention weights corresponding to context words that are most attended to for its prediction (excluding frequent words) has been highlighted with bold numerals and a light gray background. For more examples, please refer to Table 5.

Also, inspired by (Jain and Wallace, 2019), we randomly permute the attention weights and observe the change in AWE’s model performance. This idea to manipulate the attention weight distribution to study the attention mechanism dynamics is quite popular (Vashishth et al., 2019; Serrano and Smith, 2019). This experiment quantifies the change in behavior of the AWE model if it were to focus on

Sentence			professor	scolded	students	for	playing	games
Mask	professor	attention weight	-	1.112	1.666	29.096	1.125	0.171
		word vector similarity	-	0.005	0.381	-0.003	0.037	-0.228
Mask	playing	attention weight	1.392	0.932	0.641	2190.342	-	2.278
		word vector similarity	0.037	0.048	0.116	0.001	-	0.975

Sentence			mother	and	child	crossing	the	road
Mask	child	attention weight	2.527	746.334	-	0.55	0.073	0.827
		word vector similarity	1.076	0	-	-0.017	-0.005	-0.22
Mask	road	attention weight	0.353	728.901	0.553	1.86	1734.955	-
		word vector similarity	-0.05	0	-0.22	0.961	0	-

Sentence			my	elementary	school	has	a	skating	rink
Mask	skating	attention weight	0.243	0.434	1.384	39.814	274.934	-	1.354
		word vector similarity	-0.111	0.062	0.206	-0.004	-0.001	-	3.636
Mask	school	attention weight	0.138	0.972	-	21.477	1066.489	0.369	0.626
		word vector similarity	-0.011	2.595	-	-0.002	0	0.206	0.101

Table 4: Interpretation of attention. Attention weight between the masked word and the context word is given by $e^{(k_{\text{masked word}})^T q_{\text{context word}}}$, while the word vector similarity between the masked word and the context word is given by $e^{(u_{\text{masked word}})^T u_{\text{context word}}}$. Highly frequent words are highlighted with a dark gray background. For each masked word, the attention weights corresponding to context words that are most attended to (excluding highly frequent words) are highlighted with bold numerals and light gray background.

a different set of context words to predict the masked word. We find that over the entire dataset, the prediction loss for the masked word drops by 26% if we shuffle the attention weights. Thus, we can conclude that the attention mechanism in AWE focuses on informative context words for the prediction of the masked word. This is consistent with our preceding findings through the visual analysis of the attention weights (Table 4 and 5).

4.4 Role of subword information

Subwords play a critical function in the robust learning of infrequent words. For instance, the word *happiest* may not occur as frequently as the word *happy*, but the model can learn more about the word *happiest* if it were to supplement it with the information of the word *happy* as well. This is the reason why fastText and AWE-S perform significantly better than other embedding methods on the RareWords dataset (Table 1 and Table 2). FastText relates the two words *happy* and *happiest* through the n-gram ‘*happ*’. On the other hand, AWE-S relates the word *happiest* to *happy* through its verb lemma form, i.e. ‘*happy*’.

4.5 Application of AWE’s distributional loss

Deep NLP models use an input embedding matrix to project a word’s one-hot vector to a low-dimensional dense representation, which is then fed into the model. We propose to enrich this dense word vector representation with contextual information using a distributional loss term.

First, in addition to the input embedding matrix, which we refer to as the value matrix $V \in \mathbb{R}^{N \times D}$, we use two extra matrices: a key matrix $K \in \mathbb{R}^{N \times D}$ and a query matrix $Q \in \mathbb{R}^{N \times D}$. N is the size of the vocabulary and D is the embedding dimension. Let $\{w_0, w_1, \dots, w_n\}$ be the input sequence to a deep NLP model. Let w_i be an arbitrary word in the input sequence with its context given by the sequence $\{w_{i-b}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+b}\}$, where b is the size of the context window. As previously stated, for w_i , we use its value embedding v_{w_i} as input to the model. Second, we add a new loss term to the model

Sentence		colored	leaves	make	autumn	beautiful						
Mask	att. wt.	0.905	2.828	0.069	-	0.561						
	sim	-0.006	0.562	-0.116	-	0.157						
soldier	att. wt.	-	1.593	0.861	0.044	480.195	2.379					
	sim	-	0.065	-0.002	-0.108	-0.001	0.67					
oil	att. wt.	1.202	-	0.446	599.004	0.587	1.864	599.004	2.965			
	sim	0.274	-	-0.142	-0.001	-0.071	0.744	-0.001	0.973			
book	att. wt.	1.559	0.104	2090.525		0.65	2.74	31.606	4.858	0.42		
	sim	0.042	-0.292	0.001		-0.004	0.952	-0.004	0.529	0.088		
watch	att. wt.	-	1.977	32.681	0.388	3.47	1.855	714.149	1.575	159.371	32.681	1.29
	sim	-	-0.005	-0.004	-0.188	0.921	0.223	0	0.28	-0.002	-0.004	0.097
phone	att. wt.	-	5.54	755.948	1.735	2.634	176.562	0.902	1.205	567.022	755.948	4.791
	sim	-	-0.003	0	0.716	1.19	-0.002	0.002	0.508	-0.001	0	0.869

Table 5: More examples of attention weights between the masked word and the context words. Attention weight (att. wt.) between the masked word and context word is given by $e^{(\mathbf{k}_{\text{masked word}})^T \mathbf{q}_{\text{context word}}}$, while the word vector similarity (sim.) between the masked word and context word is given by $e^{(\mathbf{u}_{\text{masked word}})^T \mathbf{u}_{\text{context word}}}$. Highly frequent words are highlighted with a dark gray background. For each masked word, the attention weights corresponding to context words that are most attended to (excluding highly frequent words) are highlighted with bold numerals and a light gray background.

which enforces the input vector representation of w_i , given by \mathbf{v}_{w_i} , to be close to the weighted vector representation of its context. The loss for w_i is given by

$$l_{w_i} = 1 - \text{sigmoid}\left(\mathbf{v}_{w_i} \cdot \left(\sum_{j \in [i-b, i+b] - \{i\}} a_{w_j} \mathbf{v}_{w_j}\right)\right),$$

where,

$$a_{w_j} = \text{softmax}(\mathbf{k}_{w_i}^T \mathbf{q}_{w_j}).$$

This loss is inspired from the distributional hypothesis, which dictates that a word is known through its context (Harris, 1954). It is interesting to note that this loss can be added to any NLP model and requires no additional training data. The idea is similar to the training paradigm of word2vec and AWE, which learn word embeddings using a variant of this loss.

To test the efficacy of this loss, we train a transformer model (Vaswani et al., 2017) on WMT17 English to German news translation dataset¹. The best trained transformer model gives a perplexity score of 7.17, however, we achieve a perplexity score of 6.56 (lower is better), with the embedding scheme augmented with the new distributional loss term.

5 Conclusion & Future Work

In this work, we have proposed AWE and AWE-S, which perform significantly better than CBOW, Skip-Gram, GloVe, and fastText on a variety of datasets across a diverse set of tasks. The simple setting of masked word prediction in AWE also helped us visually analyze and interpret how the attention mechanism works. Our analysis revealed that the attention mechanism can figure out the context words that are most relevant for predicting the masked word. In future, we want to explore more applications of the proposed distributional loss for bigger models like BERT on larger datasets.

¹<http://data.statmt.org/wmt17/translation-task/>

References

- Eneko Agirre, Enrique Alfonseca, Keith Hall, Jana Kravalova, Marius Paşca, and Aitor Soroa. 2009. A study on similarity and relatedness using distributional and WordNet-based approaches. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 19–27.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015*.
- Amir Bakarov. 2018. A survey of word embeddings evaluation methods. *arXiv preprint arXiv:1801.09536*.
- 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.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642.
- Elia Bruni, Nam-Khanh Tran, and Marco Baroni. 2014. Multimodal distributional semantics. *Journal of Artificial Intelligence Research*, 49:1–47.
- Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14.
- Alexis Conneau and Douwe Kiela. 2018. SentEval: An evaluation toolkit for universal sentence representations. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc Le, and Ruslan Salakhutdinov. 2019. Transformer-XL: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2978–2988, Florence, Italy.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Bill Dolan, Chris Quirk, and Chris Brockett. 2004. Unsupervised construction of large paraphrase corpora: Exploiting massively parallel news sources. In *Proceedings of the 20th international conference on Computational Linguistics*, page 350. Association for Computational Linguistics.
- Lev Finkelstein, Evgeniy Gabrilovich, Yossi Matias, Ehud Rivlin, Zach Solan, Gadi Wolfman, and Eytan Ruppín. 2001. Placing search in context: The concept revisited. In *Proceedings of the 10th international conference on World Wide Web*, pages 406–414.
- J. R. Firth. 1957. A synopsis of linguistic theory 1930-55. *Studies in Linguistic Analysis (special volume of the Philological Society)*, 1952-59:1–32.
- Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256.
- Zellig S Harris. 1954. Distributional structure. *Word*, 10(2-3):146–162.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034.
- Felix Hill, Roi Reichart, and Anna Korhonen. 2015. Simlex-999: Evaluating semantic models with (genuine) similarity estimation. *Computational Linguistics*, 41(4):665–695.
- Sarthak Jain and Byron C. Wallace. 2019. Attention is not Explanation. 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 3543–3556.

- Stanisław Jastrzebski, Damian Leśniak, and Wojciech Marian Czarnecki. 2017. How to evaluate word embeddings? on importance of data efficiency and simple supervised tasks. *arXiv preprint arXiv:1702.02170*.
- Wang Ling, Chris Dyer, Alan W Black, and Isabel Trancoso. 2015a. Two/too simple adaptations of word2vec for syntax problems. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1299–1304.
- Wang Ling, Yulia Tsvetkov, Silvio Amir, Ramon Fernandez, Chris Dyer, Alan W Black, Isabel Trancoso, and Chu-Cheng Lin. 2015b. Not all contexts are created equal: Better word representations with variable attention. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1367–1372.
- Minh-Thang Luong, Richard Socher, and Christopher D Manning. 2013. Better word representations with recursive neural networks for morphology. In *Proceedings of the Seventeenth Conference on Computational Natural Language Learning*, pages 104–113.
- Florian Mai, Lukas Galke, and Ansgar Scherp. 2019. Cbow is not all you need: Combining cbow with the compositional matrix space model. *arXiv preprint arXiv:1902.06423*.
- Marco Marelli, Stefano Menini, Marco Baroni, Luisa Bentivogli, Raffaella Bernardi, Roberto Zamparelli, et al. 2014. A sick cure for the evaluation of compositional distributional semantic models. In *LREC*, pages 216–223.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. Efficient estimation of word representations in vector space. In Yoshua Bengio and Yann LeCun, editors, *1st International Conference on Learning Representations, ICLR 2013*.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013b. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Kira Radinsky, Eugene Agichtein, Evgeniy Gabilovich, and Shaul Markovitch. 2011. A word at a time: computing word relatedness using temporal semantic analysis. In *Proceedings of the 20th international conference on World wide web*, pages 337–346.
- Herbert Rubenstein and John B Goodenough. 1965. Contextual correlates of synonymy. *Communications of the ACM*, 8(10):627–633.
- Timo Schick and Hinrich Schütze. 2019. Attentive mimicking: Better word embeddings by attending to informative contexts. 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 489–494.
- Sofia Serrano and Noah A. Smith. 2019. Is attention interpretable? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2931–2951, July.
- Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. 2013. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pages 1139–1147.
- Shikhar Vashishth, Shyam Upadhyay, Gaurav Singh Tomar, and Manaal Faruqui. 2019. Attention interpretability across nlp tasks. *arXiv preprint arXiv:1909.11218*.
- 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.
- Wenya Wang, Sinno Jialin Pan, Daniel Dahlmeier, and Xiaokui Xiao. 2016. Recursive neural conditional random fields for aspect-based sentiment analysis. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 616–626.
- Sarah Wiegrefe and Yuval Pinter. 2019. Attention is not not explanation. 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 11–20.
- Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. 2019. Self-attention generative adversarial networks. In *Proceedings of the 36th International Conference on Machine Learning*, Proceedings of Machine Learning Research, pages 7354–7363.