

SENSE2VEC - A FAST AND ACCURATE METHOD FOR WORD SENSE DISAMBIGUATION IN NEURAL WORD EMBEDDINGS.

Andrew Trask & Phil Michalak & John Liu

Digital Reasoning Systems, Inc.

Nashville, TN 37212, USA

{andrew.trask, phil.michalak, john.liu}@digitalreasoning.com

ABSTRACT

Neural word representations have proven useful in Natural Language Processing (NLP) tasks due to their ability to efficiently model complex semantic and syntactic word relationships. However, most techniques model only one representation per word, despite the fact that a single word can have multiple meanings or "senses". Some techniques model words by using multiple vectors that are clustered based on context. However, recent neural approaches rarely focus on the application to a consuming NLP algorithm. Furthermore, the training process of recent word-sense models is expensive relative to single-sense embedding processes. This paper presents a novel approach which addresses these concerns by modeling multiple embeddings for each word based on supervised disambiguation, which provides a fast and accurate way for a consuming NLP model to select a sense-disambiguated embedding. We demonstrate that these embeddings can disambiguate both contrastive senses such as nominal and verbal senses as well as nuanced senses such as sarcasm. We further evaluate Part-of-Speech disambiguated embeddings on neural dependency parsing, yielding a greater than 8% average error reduction in unlabeled attachment scores across 6 languages.

1 INTRODUCTION

NLP systems seek to automate the extraction of information from human language. A key challenge in this task is the complexity and sparsity in natural language, which leads to a phenomenon known as the curse of dimensionality. To overcome this, recent work has learned real valued, distributed representations for words using neural networks (G.E. Hinton, 1986; Bengio et al., 2003; Morin & Bengio, 2005; Mnih & Hinton, 2009). These "neural language models" embed a vocabulary into a smaller dimensional linear space that models "the probability function for word sequences, expressed in terms of these representations" (Bengio et al., 2003). The result is a vector-space model (VSM) that represents word meanings with vectors that capture the semantic and syntactic information of words (Maas & Ng, 2010). These distributed representations model shades of meaning across their dimensions, allowing for multiple words to have multiple real-valued relationships encoded in a single vector (Liang & Potts, 2015).

Various forms of distributed representations have shown to be useful for a wide variety of NLP tasks including Part-of-Speech tagging, Named Entity Recognition, Analogy/Similarity Querying, Transliteration, and Dependency Parsing (Al-Rfou et al., 2013; Al-Rfou et al., 2015; Mikolov et al., 2013a;b; Chen & Manning, 2014). Extensive research has been done to tune these embeddings to various tasks by incorporating features such as character (compositional) information, word order information, and multi-word (phrase) information (Ling et al., 2015; Mikolov et al., 2013c; Zhang et al., 2015; Trask et al., 2015).

Despite these advancements, most word embedding techniques share a common problem in that each word must encode all of its potential meanings into a single vector (Huang et al., 2012). For words with multiple meanings (or "senses"), this creates a superposition in vector space where a vector takes on a mixture of its individual meanings. In this work, we will show that this superposition

obscures the context specific meaning of a word and can have a negative effect on NLP classifiers leveraging the superposition as input data. Furthermore, we will show that disambiguating multiple word senses into separate embeddings alleviates this problem and the corresponding confusion to an NLP model.

2 RELATED WORK

2.1 WORD2VEC

Mikolov et al. (2013a) proposed two simple methods for learning continuous word embeddings using neural networks based on Skip-gram or Continuous-Bag-of-Word (CBOW) models and named it word2vec. Word vectors built from these methods map words to points in space that effectively encode semantic and syntactic meaning despite ignoring word order information. Furthermore, the word vectors exhibited certain algebraic relations, as exemplified by example: " $v[\text{man}] - v[\text{king}] + v[\text{queen}] \approx v[\text{woman}]$ ". Subsequent work leveraging such neural word embeddings has proven to be effective on a variety of natural language modeling tasks (Al-Rfou et al., 2013; Al-Rfou et al., 2015; Chen & Manning, 2014).

2.2 WANG2VEC

Because word embeddings in word2vec are insensitive to word order, they are suboptimal when used for syntactic tasks like POS tagging or dependency parsing. Ling et al. (2015) proposed modifications to word2vec that incorporated word order. Consisting of structured skip-gram and continuous window methods that are together termed wang2vec, these models demonstrate significant ability to model syntactic representations. They come, however, at the cost of computation speed. Furthermore, because words have a single vector representation in wang2vec, the method is unable to model polysemic words with multiple meanings. For instance, the word "work" in the sentence "We saw her work" can be either a verb or noun depending on the broader context in surrounding this sentence. This technique encodes the co-occurrence statistics for each sense of a word into one or more fixed dimensional embeddings, generating embeddings that model multiple uses of a word.

2.3 STATISTICAL MULTI-PROTOTYPE VECTOR-SPACE MODELS OF WORD MEANING

Perhaps a seminal work to vector-space word-sense disambiguation, the approach by Reisinger & Mooney (2010) creates a vector-space model that encodes multiple meanings for words by first clustering the contexts in which a word appears. Once the contexts are clustered, several prototype vectors can be initialized by averaging the statistically generated vectors for each word in the cluster. This process of computing clusters and creating embeddings based on a vector for each cluster has become the canonical strategy for word-sense disambiguation in vector spaces. However, this approach presents no strategy for the context specific selection of potentially many vectors for use in an NLP classifier.

2.4 CLUSTERING WEIGHTED AVERAGE CONTEXT EMBEDDINGS

Our technique is inspired by the work of Huang et al. (2012), which uses a multi-prototype neural vector-space model that clusters contexts to generate prototypes. Unlike Reisinger & Mooney (2010), the context embeddings are generated by a neural network in the following way: given a pre-trained word embedding model, each context embedding is generated by computing a weighted sum of the words in the context (weighted by tf-idf). Then, for each term, the associated context embeddings are clustered. The clusters are used to re-label each occurrence of each word in the corpus. Once these terms have been re-labeled with the cluster's number, a new word model is trained on the labeled embeddings (with a different vector for each) generating the word-sense embeddings.

In addition to the selection problem and clustering overhead described in the previous subsection, this model also suffers from the need to train neural word embeddings twice, which is a very expensive endeavor.

2.5 CLUSTERING CONVOLUTIONAL CONTEXT EMBEDDINGS

Recent work has explored leveraging convolutional approaches to modeling the context embeddings that are clustered into word prototypes. Unlike previous approaches, Chen et al. (2015) selects the number of word clusters for each word based on the number of definitions for a word in the WordNet Gloss (as opposed to other approaches that commonly pick a fixed number of clusters). A variant on the MSSG model of Neelakantan et al. (2015), this work uses the WordNet Glosses dataset and convolutional embeddings to initialize the word prototypes.

In addition to the selection problem, clustering overhead, and the need to train neural embeddings multiple times, this higher-quality model is somewhat limited by the vocabulary present in the English WordNet resource. Furthermore, the majority of the WordNets relations connect words from the same Part-of-Speech (POS). "Thus, WordNet really consists of four sub-nets, one each for nouns, verbs, adjectives and adverbs, with few cross-POS pointers."¹

3 THE SENSE2VEC MODEL

We expand on the work of Huang et al. (2012) by leveraging supervised NLP labels instead of unsupervised clusters to determine a particular word instance's sense. This eliminates the need to train embeddings multiple times, eliminates the need for a clustering step, and creates an efficient method by which a supervised classifier may consume the appropriate word-sense embedding.

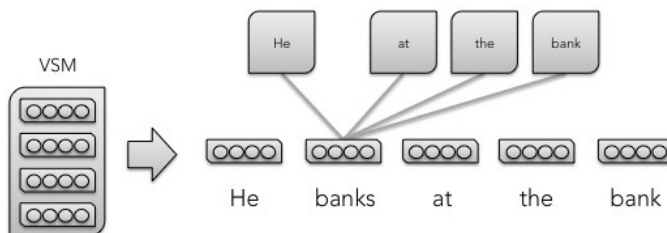


Figure 1: A graphical representation of wang2vec.

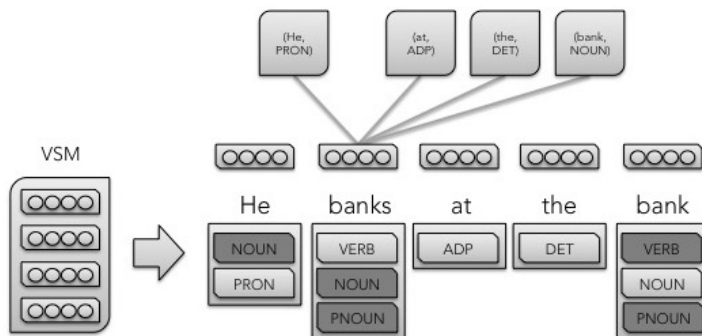


Figure 2: A graphical representation of sense2vec.

Given a labeled corpus (either by hand or by a model) with one or more labels per word, the sense2vec model first counts the number of uses (where a unique word maps set of one or more

¹<https://wordnet.princeton.edu/>

labels/uses) of each word and generates a random "sense embedding" for each use. A model is then trained using either the CBOW, Skip-gram, or Structured Skip-gram model configurations. Instead of predicting a token given surrounding tokens, this model predicts a word sense given surrounding senses.

3.1 SUBJECTIVE EVALUATION - SUBJECTIVE BASELINE

For subjective evaluation of these word embeddings, we trained models using several datasets for comparison. First, we trained using Word2vec's Continuous Bag of Words² approach on the large unlabeled corpus used for the Google Word Analogy Task³. Several word embeddings and their closest terms measured by cosine similarity are displayed in Table 1 below.

Table 1: Single-sense Baseline Cosine Similarities

bank	1.0	apple	1.0	so	1.0	bad	1.0	perfect	1.0
banks	.718	iphone	.687	but	.879	good	.727	perfection	.681
banking	.672	ipad	.649	it	.858	worse	.718	perfectly	.670
hsbc	.599	microsoft	.603	if	.842	lousy	.717	ideal	.644
citibank	.586	ipod	.595	even	.833	stupid	.710	flawless	.637
lender	.566	imac	.594	do	.831	horrible	.703	good	.622
lending	.559	iphones	.578	just	.808	awful	.697	always	.572

In this table, observe that the "bank" column is similar to proper nouns ("hsbc", "citibank"), verbs ("lending", "banking"), and nouns ("banks", "lender"). This is because the term "bank" is used in 3 different ways, as a proper noun, verb, and noun. This embedding for "bank" has modeled a mixture of these three meanings. "apple", "so", "bad", and "perfect" can also have a mixture of meanings. In some cases, such as "apple", one interpretation of the word is completely ignored (apple the fruit). In the case of "so", there is also an interjection sense of "so" that is not well represented in the vector space.

3.2 SUBJECTIVE EVALUATION - PART-OF-SPEECH DISAMBIGUATION

For Part-of-Speech disambiguation, we labeled the dataset from section 3.1 with Part-of-Speech tags using the Polyglot Universal Dependency Part-of-Speech tagger of Al-Rfou et al. (2013) and trained sense2vec with identical parameters as section 3.1. In table 2, we see that this method has successfully disambiguated the difference between the noun "apple" referring to the fruit and the proper noun "apple" referring to the company. In table 3, we see that all three uses of the word "bank" have been disambiguated by their respective parts of speech, and in table 4, nuanced senses of the word "so" have also been disambiguated.

Table 2: Part-of-Speech Cosine Similarities for the Word: apple

apple	NOUN	1.0	apple	PROPN	1.0
apples	NOUN	.639	microsoft	PROPN	.603
pear	NOUN	.581	iphone	NOUN	.591
peach	NOUN	.579	ipad	NOUN	.586
blueberry	NOUN	.570	samsung	PROPN	.572
almond	NOUN	.541	blackberry	PROPN	.564

²command line params: -size 500 -window 10 -negative 10 -hs 0 -sample 1e-5 -iter 3 -min-count 10

³the data.txt file generated from <http://word2vec.googlecode.com/svn/trunk/demo-train-big-model-v1.sh>

Table 3: Part-of-Speech Cosine Similarities for the Word: bank

bank	NOUN	1.0	bank	PROPN	1.0	bank	VERB	1.0
banks	NOUN	.786	bank	NOUN	.570	gamble	VERB	.533
banking	NOUN	.629	hsbc	PROPN	.536	earn	VERB	.485
lender	NOUN	.619	citibank	PROPN	.523	invest	VERB	.470
bank	PROPN	.570	wachovia	PROPN	.503	reinvest	VERB	.466
ubs	PROPN	.535	grindlays	PROPN	.492	donate	VERB	.466

Table 4: Part-of-Speech Cosine Similarities for the Word: so

so	INTJ	1.0	so	ADV	1.0	so	ADJ	1.0
now	INTJ	.527	too	ADV	.753	poved	ADJ	.588
obviously	INTJ	.520	but	CONJ	.752	condemnable	ADJ	.584
basically	INTJ	.513	because	SCONJ	.720	disputable	ADJ	.578
okay	INTJ	.505	but	ADV	.694	disapprove	ADJ	.559
actually	INTJ	.503	really	ADV	.671	contestable	ADJ	.558

3.3 SUBJECTIVE EVALUATION - SENTIMENT DISAMBIGUATION

For Sentiment disambiguation, the IMDB labeled training corpus was labeled with Part-of-Speech tags using the Polyglot Part-of-Speech tagger from Al-Rfou et al. (2013). Adjectives were then labeled with the positive or negative sentiment associated with each comment. A CBOW sense2vec model was then trained on the resulting dataset, disambiguating between both Part-of-Speech and Sentiment (for adjectives).

Table 5 shows the difference between the positive and negative vectors for the word "bad". The negative vector is most similar to word indicating the classical meaning of bad (including the negative version of "good", e.g. "good grief!"). The positive "bad" vector denotes a tone of sarcasm, most closely relating to the positive sense of "good" (e.g. "good job!").

Table 5: Sentiment Cosine Similarities for the Word: bad

bad	NEG	1.0	bad	POS	1.0
terrible	NEG	.905	good	POS	.753
horrible	NEG	.872	wrong	POS	.752
awful	NEG	.870	funny	POS	.720
good	NEG	.863	great	POS	.694
stupid	NEG	.845	weird	POS	.671

Table 6 shows the positive and negative senses of the word "perfect". The positive version of the word clusters most closely with words indicating excellence. The positive version clusters with the more sarcastic interpretation.

Table 6: Sentiment Cosine Similarities for the Word: perfect

perfect	NEG	1.0	perfect	POS	1.0
real	NEG	0.682	wonderful	POS	0.843
unfortunate	NEG	0.680	brilliant	POS	0.842
serious	NEG	0.673	incredible	POS	0.840
complete	NEG	0.673	fantastic	POS	0.839
ordinary	NEG	0.673	great	POS	0.823
typical	NEG	0.661	excellent	POS	0.822
misguided	NEG	0.650	amazing	POS	0.814

4 NAMED ENTITY RESOLUTION

To evaluate the embeddings when disambiguating on named entity resolution (NER), we labeled the standard word2vec dataset from section 3.2 with named entity labels. This demonstrated how sense2vec can also disambiguate between multi-word sequences of text as well as single word sequences of text. Below, we see that the word "Washington" is disambiguated with both a PERSON and a GPE sense of the word. Furthermore, we see that Hillary Clinton is very similar to titles that she has held within the time span of the dataset.

Table 7: Disambiguation for the word: Washington

George_Washington	PERSON_NAME	.656	Washington_D	GPE	.665
Henry_Knox	PERSON_NAME	.624	Washington_DC	GPE	.591
Philip_Schuyler	PERSON_NAME	.618	Seattle	GPE	.559
Nathanael_Greene	PERSON_NAME	.613	Warsaw_Embassy	GPE	.524
Benjamin_Lincoln	PERSON_NAME	.602	Wash	GPE	.516
William_Howe	PERSON_NAME	.591	Maryland	GPE	.507

Table 8: Entity resolution for the term: Hillary Clinton

Secretary_of_State	TITLE	0.661
Senator	TITLE	0.613
Senate	ORG_NAME	0.564
Chief	TITLE	0.555
White_House	ORG_NAME	0.564
Congress	ORG_NAME	0.547

5 NEURAL DEPENDENCY PARSING

To quantitatively evaluate disambiguated sense embeddings relative to the current standard, we compared sense2vec embeddings and wang2vec embeddings on neural syntactic dependency parsing tasks in six languages. First, we trained two sets of embeddings on the Bulgarian, German, English, French, Italian, and Swedish Wikipedia datasets from the Polyglot website⁴. The baseline embeddings were trained without any Part-of-Speech disambiguation using the structured skip-gram approach of Ling et al. (2015). For each language, the sense2vec embeddings were trained by disambiguating terms using the language specific Polyglot Part-of-Speech tagger of Al-Rfou et al. (2013), and embedded in the same structured skip-gram approach. Both were trained using identical parametrization⁵.

⁴<https://sites.google.com/site/rmyeid/projects/polyglot>

⁵command line params: -size 50 -window 5 -negative 10 -hs 0 -sample 1e-4 -iter 5 -cap 0

Each of these embeddings was used to train a dependency parse model using the parser outlined in (Chen & Manning, 2014). All were trained on the the respective language’s Universal Dependencies treebank. The standard splits were used.⁶ For the parser trained on the sense2vec emeddings, the POS specific embedding was used as the input. The Part-of-Speech label was determined using the gold-standard POS tags from the treebank. It should be noted that the parser of (Chen & Manning, 2014) uses trained Part-of-Speech embeddings as input which are indexed based on gold-standard POS tags. Thus, differences in quality between parsers trained on the two embedding styles are due to clarity in the word embeddings as opposed to the addition of Part-of-Speech information because both model styles train on gold standard POS information. For each language, the Unlabeled Attachment Scores are outlined in Table 7.

Table 9: Unlabeled Attachment Scores and Percent Error Reductions

	Set	Bulgarian	German	English	French	Italian	Swedish	Mean
wang	Dev	90.03	68.86	85.02	73.82	84.99	78.94	80.28
	Test*	90.17	60.25	83.61	70.10	84.99	82.47	78.60
	Test	90.39	60.54	83.88	70.53	85.45	82.51	78.88
sense	Dev	90.69	72.61	86.10	75.43	85.57	81.21	81.94
	Test*	90.41	64.17	85.48	71.66	86.13	84.44	80.38
	Test	90.86	64.43	85.93	72.16	86.18	84.60	80.69
Error Margin	Dev	7.05%	13.69%	7.76%	6.56%	3.98%	12.06%	8.52%
	Test	2.47%	10.95%	12.82%	5.50%	8.21%	12.71%	8.78%
	Abs.	5.17%	10.93%	14.54%	5.86%	5.32%	13.58%	9.23%
	Avg.	4.76%	12.32%	10.29%	6.03%	6.09%	12.39%	

The ”Error Margin” section of table 7 describes the percentage reduction in error for each language. Disambiguating based on Part-of-Speech using sense2vec reduced the error in all six languages with an average reduction greater than 8%.

6 CONCLUSION AND FUTURE WORK

In this work, we have proposed a new model for word sense disambiguation that uses supervised NLP labeling to disambiguate between word senses. Much like previous models, it leverages a form of context clustering to disambiguate the use of a term. However, instead of using unsupervised clustering methods, our approach clusters using supervised labels which can analyze a specific word’s context and assign a label. This significantly reduces the computational overhead of word-sense modeling and provides a natural mechanism for other NLP tasks to select the appropriate sense embedding. Furthermore, we show that disambiguated embeddings can increase the accuracy of syntactic dependency parsing in a variety of languages. Future work will explore how disambiguated embeddings perform using other varieties of supervised labels and consuming NLP tasks.

REFERENCES

- Al-Rfou, Rami, Perozzi, Bryan, and Skiena, Steven. Polyglot: Distributed word representations for multilingual NLP. *CoRR*, abs/1307.1662, 2013. URL <http://arxiv.org/abs/1307.1662>.
- Al-Rfou, Rami, Kulkarni, Vivek, Perozzi, Bryan, and Skiena, Steven. Polyglot-NER: Massive multilingual named entity recognition. *Proceedings of the 2015 SIAM International Conference on Data Mining, Vancouver, British Columbia, Canada, April 30 - May 2, 2015*, April 2015.
- Bengio, Yoshua, Ducharme, Réjean, Vincent, Pascal, and Janvin, Christian. A neural probabilistic language model. *J. Mach. Learn. Res.*, 3:1137–1155, March 2003. ISSN 1532-4435.

⁶The German, French, and Italian treebanks had occasional tokens that both spanned multiple indices and overlapped with the index of the previous and following token (ex. 0, 0-1, 1,...), a property which is incompatible with the (Chen & Manning, 2014) parser. These tokens were removed. If their removal created a malformed tree, the sentence was removed automatically by the parser and logged accordingly.

- Chen, Danqi and Manning, Christopher. A fast and accurate dependency parser using neural networks. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 740–750, Doha, Qatar, October 2014. Association for Computational Linguistics. URL <http://www.aclweb.org/anthology/D14-1082>.
- Chen, Tao, Xu, Ruifeng, He, Yulan, and Wang, Xuan. Improving distributed representation of word sense via wordnet gloss composition and context clustering. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pp. 15–20, Beijing, China, July 2015. Association for Computational Linguistics. URL <http://www.aclweb.org/anthology/P15-2003>.
- G.E. Hinton, J.L. McClelland, D.E. Rumelhart. Distributed representations. *Parallel distributed processing: Explorations in the microstructure of cognition*, 1(3):77–109, 1986.
- Huang, Eric H., Socher, Richard, Manning, Christopher D., and Ng, Andrew Y. Improving word representations via global context and multiple word prototypes. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers - Volume 1, ACL '12*, pp. 873–882, Stroudsburg, PA, USA, 2012. Association for Computational Linguistics. URL <http://dl.acm.org/citation.cfm?id=2390524.2390645>.
- Liang, P. and Potts, C. Bringing machine learning and compositional semantics together. *Annual Reviews of Linguistics*, 1(1):355–376, 2015.
- Ling, Wang, Dyer, Chris, Black, Alan W, and Trancoso, Isabel. 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*, pp. 1299–1304, Denver, Colorado, May–June 2015. Association for Computational Linguistics. URL <http://www.aclweb.org/anthology/N15-1142>.
- Maas, Andrew L and Ng, Andrew Y. A probabilistic model for semantic word vectors. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2010.
- Mikolov, Tomas, Chen, Kai, Corrado, Greg, and Dean, Jeffrey. Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781, 2013a. URL <http://arxiv.org/abs/1301.3781>.
- Mikolov, Tomas, Le, Quoc V., and Sutskever, Ilya. Exploiting similarities among languages for machine translation. *CoRR*, abs/1309.4168, 2013b. URL <http://arxiv.org/abs/1309.4168>.
- Mikolov, Tomas, Sutskever, Ilya, Chen, Kai, Corrado, Greg, and Dean, Jeffrey. Distributed representations of words and phrases and their compositionality. *CoRR*, abs/1310.4546, 2013c. URL <http://arxiv.org/abs/1310.4546>.
- Mnih, Andriy and Hinton, Geoffrey E. A scalable hierarchical distributed language model. In Koller, D., Schuurmans, D., Bengio, Y., and Bottou, L. (eds.), *Advances in Neural Information Processing Systems 21*, pp. 1081–1088. Curran Associates, Inc., 2009.
- Morin, Frederic and Bengio, Yoshua. Hierarchical probabilistic neural network language model. In *Proceedings of the international workshop on artificial intelligence and statistics*, pp. 246–252. Citeseer, 2005.
- Neelakantan, Arvind, Shankar, Jeevan, Passos, Alexandre, and McCallum, Andrew. Efficient non-parametric estimation of multiple embeddings per word in vector space. *CoRR*, abs/1504.06654, 2015. URL <http://arxiv.org/abs/1504.06654>.
- Reisinger, Joseph and Mooney, Raymond J. Multi-prototype vector-space models of word meaning. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics, HLT '10*, pp. 109–117, Stroudsburg, PA, USA, 2010. Association for Computational Linguistics. ISBN 1-932432-65-5. URL <http://dl.acm.org/citation.cfm?id=1857999.1858012>.

Trask, Andrew, Gilmore, David, and Russell, Matthew. Modeling order in neural word embeddings at scale. *CoRR*, abs/1506.02338, 2015. URL <http://arxiv.org/abs/1506.02338>.

Zhang, Xiang, Zhao, Junbo, and LeCun, Yann. Character-level convolutional networks for text classification. *CoRR*, abs/1509.01626, 2015. URL <http://arxiv.org/abs/1509.01626>.