

Evaluation of Word Embeddings for Sequence Tagging Tasks

A Anonymous

B Anonymous

Abstract

Word embeddings algorithms have the ability to learn word representation from unlabelled data. The resulting representations are used as features in many NLP applications. In this paper, we carry on an extrinsic evaluation of five different word embedding methods under control conditions, and evaluate them in four sequence labelling tasks: POS-tagging, chunking, NER and MWE identification. We also evaluate the performance of fine-tuned versus non fine-tuned features during training, and show that fine-tuning can result in over-fitting, when the representations learned unsupervised were already good. We also found that when using word embeddings, only several hundred of training instances are needed to reach a decent performance. Surprisingly, we could not find any leading word embedding method across the different tasks and proposed settings. Nevertheless, word embeddings are always helping to improve the performance and are especially useful for task or domains with limited training data. They are also easy to train and to integrate to existing NLP systems, and therefore, should be considered in modern NLP architectures.

1 Introduction

Tim

In the last years, distributed word representations have been applied to several NLP tasks. Inspired by distributional semantics models, distributed word representation methods represent each word as a continuous vector, where similar words have a similar vector representation, therefore, capturing the similarity between words.

The resulting vectors can be used as features in many NLP applications and it has been shown that they outperform methods that treat words as atomic units (\emptyset). Their attractiveness relies in the ability to learn word representations in an unsupervised way, thus directly providing lexical features from big amounts of unlabelled data and, therefore, alleviating the cost of human annotations. It has been also claimed that word embeddings have the ability to connect out of vocabulary words to known ones. Hence, suggesting that word embeddings are a good resource for applications that need to be adapted to a certain domain, different from the one the application has been tuned for. For example,... Another property attribute to word embeddings is their capacity to encourage common behaviour among related in-vocabulary words, for instance...

As with other learning methods, it is well known that the performance of machine learning algorithms heavily depends on parameter optimization, the size of the training data used and the applications they target. For example, (Turian et al., 2010) shows that the optimal word embedding dimensions are task specific. Moreover, there are several word embeddings methods, which used different algorithms and resources. Some methods involve feedback from the end task when learning (or fine-tuning) the word representations and others do not. Learning algorithms that involve fine-tuning are supposed to perform better since word representations become task-specific, at the cost of performing worst for out of vocabulary words. But still, there is not systematic comparison between these two methods.

In this paper, we perform an extensive evaluation of five word embedding approaches under fixed experiment conditions, and evaluate them over different sequence labelling tasks: POS-tagging, chunking, NER and MWE (Multi Word Expression Identification), within the following

aims: (i) perform a fair comparison of different word embeddings algorithms. This includes running different word embeddings algorithms under controlled conditions, for example, use the same training set, the same preprocessing, etc.; (ii) measure the influence of word embeddings in sequence labeling tasks in semi-supervised settings (fine-tuning); (iii) systematically compared the usefulness of word embedding versus unigram features for sequence tagging. (iv) use word embeddings for MWE. To the best of our knowledge, word embeddings have not been used for this task before;

2 Self-taught Learning for Sequence Tagging

This section might go to introduction The idea of learning word representations for downstream NLP applications embraces the idea of self-taught learning (). Self-taught learning starts with learning initial data representations from a large amount of unlabelled data, which may not be directly relevant to target applications. The learned representations are then fed as features into models of target applications, and may be fine-tuned during training to adapt to the target needs. Learning from a random sample of unlabelled data is distinct from semi-supervised learning (), which makes an assumption that the unlabelled data shares the same distribution as the labelled training data. In the following, we will introduce the recent advances of learning word representations and apply them for a range of sequence tagging tasks by using graph transformers ().

3 Word Representations

The distributional hypothesis in linguistics suggests that "a word is characterised by the company it keeps" (Firth, 1957). Words that are used in the similar contexts tend to have similar semantic and syntactic properties. Capturing distributional similarity is the underlying idea of all word representation learning methods.

3.1 Types of Word Representations

Based on the ways of constructing word representations, Turian et al. (2010) categorise these methods into three types: *Distributional representation*, *Cluster-based representation*, and *Distributed representation*.

Distributional representation methods map each word w to its context word vector C_w , which is built based on cooccurrence counts between w and the words surrounding it. The learning methods store either directly the cooccurrence count between two words w and i in C_{wi} () or project the concurrence counts between words into a lower dimensional space by using techniques such as SVD () and LDA ().

The methods of *Cluster-based representation* build clusters of words by applying either soft- or hard clustering algorithms. Some of them also rely on cooccurrence matrix of words (). The Brown algorithm () is the most famous one in this category.

A *distributed representation* takes the form of a dense, low-dimensional, and continuous-valued vector. It is compact and stores latent features of a word. This kind of representations are learned with various neural language models () in the hope of capturing both syntactic and semantic properties of words.

3.2 Selected Word Representations

For the sequence tagging tasks, we choose five top performed word representations in various empirical studies as candidates : Brown clustering (), CW (), Skip-gram (), CBOW (), and Glove (). The key idea of all these word models is to estimate the probability of word sequences. Formally, given a word sequence $\mathbf{w} = \langle w_1, w_2, \dots, w_T \rangle$, the training objective of these models is to maximise the log probability.

$$p(\mathbf{w}) = \frac{1}{T} \sum_{k=1}^T \log(p(w_k | w_{j \in c_{k-m}^{k+n}})) \quad (1)$$

where c_{k-m}^{k+n} denotes a sub-sequence $\langle w_{k-m}, \dots, w_{k-1}, w_k, w_{k+1}, \dots, w_{k+n} \rangle$ of length $n - m - 1$, which is the local context of w_k . All models except CW choose $n = -1$ and $m > 0$ and exclude the word w_k from the local context.

The key differences among these models are the parameterisation of the factors $p(w_k | w_{j \in c_{k-m}^{k+n}})$ as well as the training loss functions. Brown clustering introduces a finite set of word classes V for each word,

3.3 Building Word Representations

For a fair comparison, we train each kind of word embedding on a combination of all corpora in Table 1. The joint corpus was preprocessed with the

Stanford sentence splitter and tokenizer. All consecutive digit substrings were replaced by NUM_f , where f is the length of the digit substring (e.g., “10.20” is replaced by “ $\text{NUM}_2.\text{NUM}_2$ ”).

Data set	Size	Words
UMBC	48.1GB	3 billions
One Billion	4.1GB	0.8 billions
Latest wikipedia dump	49.6GB	3 billions

Table 1: Corpora used to learn word embeddings

4 Sequence Tagging Tasks

Lizhen

We evaluate different word representations in four different sequence tagging tasks: POS tagging, chunking, NER and MWE identification.

For each sequence tagging task, we feed learned word representations into a first order linear-chain graph transformer (Collobert et al., 2011), and trained them by using the online learning algorithm AdaGrad (Duchi et al., 2011). For each model taking distributed word representations as word features, we consider two settings:

- Graph transformer *does not* fine tune word representations during training (this is equivalent to a linear-chain CRF);
- Graph transformer fine tunes the word representations during training.

We consider also CRF models with hand-crafted features, which use one-hot representation for each unigram.

We split the task specific corpus into a training set, validation set, and a test set (see Table 2). If a corpus already provides fixed splits, we reuse them. For POS-tagging and NER, we also evaluated the models with a out-of-domain corpus (English Web-Treebank and MUC-7, respectively), which has similar annotation schema as the respective training corpus.

In order to have fair and reproducible experimental results, we tuned the hyperparameters with random search (Bergstra and Bengio, 2012). We randomly sampled 50 distinct hyperparameter sets with the same random seed for the models that do not update word embeddings, and sampled 100 distinct hyperparameter sets for the models that update word embeddings. For each set of hyperparameters, we train a model on its training set and pick up the best one based on its performance on its validation set (Turian et al., 2010). Note that,

we also consider word vector size and context window size of distributed word representation, and number of clusters of brown clustering as the hyperparameters. This is achieved by mapping each possible hyperparameter combination to the word representation files trained with these parameters. Their ranges are listed below.

- **Word vector size:** [25, 50, 100, 200].
- **Context window size:** [5, 10, 15].
- **Number of brown clusters:** [250, 500, 1000, 2000, 4000].

However, for the models that update word representations, we always found under-performed hyperparameters after trying out all hyperparameter combinations, because they have more hyperparameters than the models that do not update word representations. Then, for each distributed word representations, we reuse all hyperparameters of the models that do not update word representations, only tune the hyperparameters of AdaGrad for the word representation layer. This method requires only 32 additional runs for each model updating embeddings and achieves consistently better results than 100 random draws.

The final evaluation is carried out in a semi-supervised setting. We split the training set into 10 partitions at log scale. That means, the second smallest partition will be twice the size of the smallest partition. We created 10 training sets with incremental size by merging these partitions from the smallest one to the largest one, and each of them on the same designated test sets.

We adopted the most commonly used F1 measure as the evaluation metric for all tasks except POS tagging, for which we use per-word accuracy. In order to evaluate model performance on out-of-vocabulary (unknown) words, we reported also the accuracy for the words that do not occur in the training set.

In addition, we also set up experiments to verify if CRF/graph transformer requires different feature design for different kinds of pre-trained word embeddings. This is achieved by adding a hidden layer between CRF and distributed word representations. For each context word, the hidden layer computes the element-wise multiplication of its embedding with the embedding of the current word embedding, and the representation of current word stays the same. The results of this approach are not plotted because this method leads only to marginal improvement.

Table 2: Datasets splits and feature space for each sequence tagging task.

	Training set	Validation set	Test set	Feature space
POS-Tagging	0-18 WSJ	19-21 WSJ	22-24 of WSJ; English Web-Treebank	as in (Collobert et al., 2011)
Chunking	WSJ	1000 sentences WSJ	CoNLL-2000	as in (Turian et al., 2010)
NER	CoNLL-2003 train set	CoNLL-2003 dev. set	CoNLL-2003 test set; MUC7	as in (Turian et al., 2010)
MWE	500 documents from	100 documents from	123 documents	as in (Schneider et al., 2014b)

5 Experimental Results and Discussion

As a reference, we compared our best results for each task with their corresponding benchmarks (Table 3). For POS-tagging and chunking, we reach a comparable performance to the state-of-the-art methods. The difference between our NER system and its baseline is most obvious, as we are 0.025 points below them, but the comparison is not fair considering that their algorithm is much complex (Ando and Zhang (2005) used a 2nd order CRF, while we used a 1st order CRF). For the task of MWE identification, our implementation and settings beat the baseline. However, in this paper, we do not aim to maximize the absolute performance of the tasks under study, but rather to study the impact of word embeddings for sequence tagging tasks under control settings. Accordingly, we focus in addressing the following questions:

(i) Are the evaluated word embedding methods better than unigram features? To answer this question, we systematically compared the usefulness of word embedding versus unigram features for sequence tagging and noted that word embedding methods always outperform unigram features (Figures 1, 2).

(ii) How does the size of labelled training data affect the experimental results? We observed that embedding methods are especially helping POS-tagging and chunking when there are only several hundreds of training instances. Therefore, confirming that any of the evaluated word embedding methods should be use when labelled data is limited. These early improvements are less evident for NER and MWE. We attribute this to: a) NER and MWE are more difficult tasks than POS-tagging and chunking; b) NER and MWE require more training data to reach a decent performance; c) the performance of NER and MWE heavily dependent on complex features such as gazetteers and lexicons like Wordnet, which are not captured by the feature representation learned from unlabelled data, meanwhile the standard features used

in POS-tagging and chunking (e.g., context words, stemmed words) are similar to the representations encoded in the word embeddings.

(iii) How well do the word embeddings perform for unknown words? As already mentioned, we measure the performance for out-of-vocabulary words (OOV) in two settings: with *in-domain* and *out-of-domain* corpora, for all the tasks, except MWE identification for which there is no other data set available (see Table 2). As expected, word embeddings and Brown clustering excel in *out-of-domain* performance. Word embeddings without fine-tuning enhance even more the performance of OOV for the *in-domain* and *out-of-domain* settings (Figure 3) since fine-tuned word representations become task-specific, hence performing worst for OOV.

(iv) How well do different word embeddings perform in all tasks when semi-supervised fine-tuning is not performed?, and **(v) and how well do different word embeddings perform in all tasks when semi-supervised fine-tuning is performed?** Across all the methods, fine-tuning is helping POS-Tagging and MWE, where the CW method has been found to be the most sensible to tuning, reaching almost 3 points more, when tuning is performed. For chunking and NER, the best results are fine-tuned, but the difference across all the methods and updated features versus not-updated ones, is not significant. We also found that fine-tuning can correct poorly learned word representations but can be overfitted if unsupervisedly learned ones are already good.

Finally, we address the following question: **(vi) It has been shown that Brown clusters are useful features for MWE identification but, are also word embeddings helping MWE identification?** According to our experiments, the word embedding features distilled by fine-tuned CW reached the best results, beating the state-of-the-art performance (see Table 3). However, between Brown clusters and fine-tuned CW, learned under the same settings, the difference is not impressive,

Table 3: Benchmark results vs. our best results

Task	Benchmark	Us
POS-Tagging	(Accuracy) 97.24 (Toutanova et al., 2003)	0.9592 (skip-gram negsam+up)
Chunking	(F1) 0.9429 (Sha and Pereira, 2003)	0.9386 (Brown cluster v2000+)
NER	(F1) 0.8931 (Ando and Zhang, 2005)	0.8686 (skip-gram negsam+noup)
MWE	(F1) 0.6253 (Schneider et al., 2014a)	0.6546 (cw+up)

suggesting that distributional word representations and cluster-based representations captures similar features for the MWE identification task. Thus, a natural question: would it be better to learn distributional representations for MWE, instead of representations of single words?

6 Related Work

Word embedding learning methods have been applied to several NLP tasks that we summaries in this section.

Collobert et al. (2011) proposed a neuronal network architecture that learn word embeddings and use them in POS-tagging, chunking, NER and Semantic Role Labelling. Without specializing their architecture for the mentioned tasks, they achieve close state-of-the-art performance. After including specialized features (e.g., word suffixes for POS-tagging; Gazzeters for NER, etc.) and other tricks like cascading and ensambling classifiers, achieve competitive state-of-the-art performance. Similarly, Turian et al. (2010) explored the impact of using word features learned from cluster-based and word embeddings representations for NER and chunking. They conclude that unsupervised word representation improve NER and chunking, and that combining different word representations can further improve the performance. Word representation from Brown clusters have been also shown to enhance Twitter POS tagging Owoputi et al. (2013).

Schneider et al. (2014a) presented a MWE analyser that, among other features, used unsupervised word clusters. They observed that the clusters were useful for identifying words that usually belong to proper names, which are considered MWE in the data set used. Nevertheless, they mentioned that it is difficult to measure the impact of the word embeddings features, since other features may capture the same information.

Word embeddings have been also used as features for syntactic dependency parsing and constituent parsing. Bansal et al. (2014) used word

embeddings as features for dependency parsing, which used the syntactic dependency context instead of the linear context in raw text. They found that simple attempts based on discretization of individual word vector dimensions do not improve parsing. Only when performing hierarchical clustering of the continuous word vectors then using features based on the hierarchy, they gain performance. They also pointed out that ensemble of different word embeddings representations improved performance. Within the same aim, Andreas and Klein (2014) explores the used of word embeddings for constituency parsing and conclude that the information they provide might be redundant with the one acquire by a syntactic parser trained with a small amount of data. Others that boost the performance when including word embeddings representations for syntactic parsing includes (Koo et al., 2008; Koo et al., 2010; Haffari et al., 2011; Tratz and Hovy, 2011).

Word embedding have also been applied to other (non-sequential NLP) tasks such as super-sense tagging (Edouard Grave and Bach, 2013); grammar induction (Spitkovsky et al., 2011) and semantic task such as semantic relatedness, synonymy detection, concept categorization selection preferences and analogy (Baroni et al., 2014)

7 Conclusion

We have performed an extensive extrinsic evaluation of five word embeddings methods approaches under fixed experiments conditions, and evaluate them over different sequence tagging tasks: POS-tagging, chunking, NER and MWE identification. We found that word embeddings methods always outperformed unigram features, especially when the training size is small, but no method was consistently better than the others across the different tasks and settings. Word representations were also found useful for improving the accuracy of OOV. We expected to see an important gap between the performance of fine-tuned features in a semi-supervised setting and no fine-tuned ones,

Figure 1: Best results for each method for POS-Tagging and Chunking. The x-axis correspond to the different word embeddings methods and the y-axis to the 10 training partitions at log scale. Green color stand for high performance, while red color stands for low performance. The methods are in chronological order

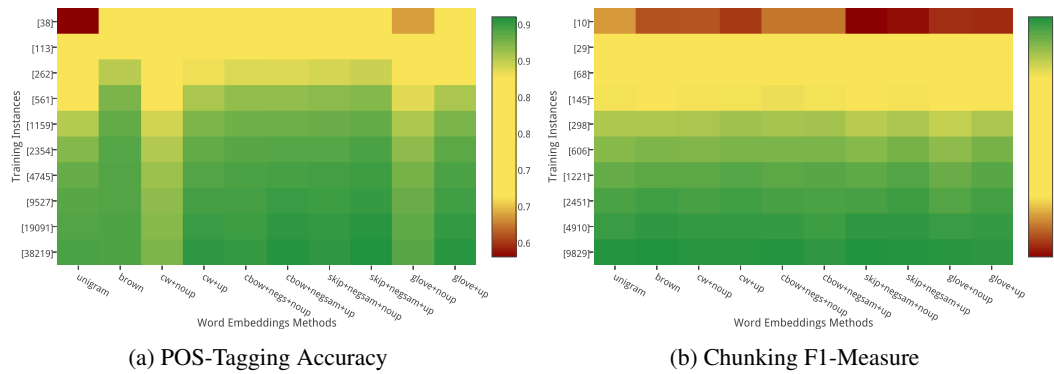


Figure 2: Best results for each method for NER and MWE

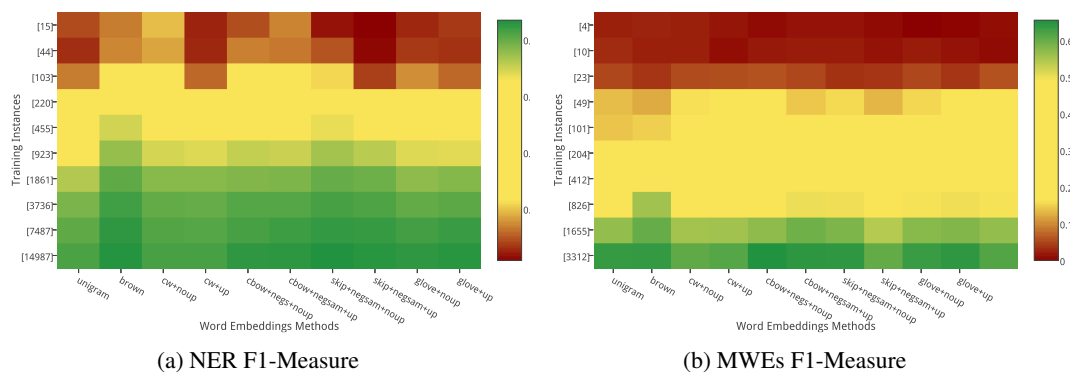
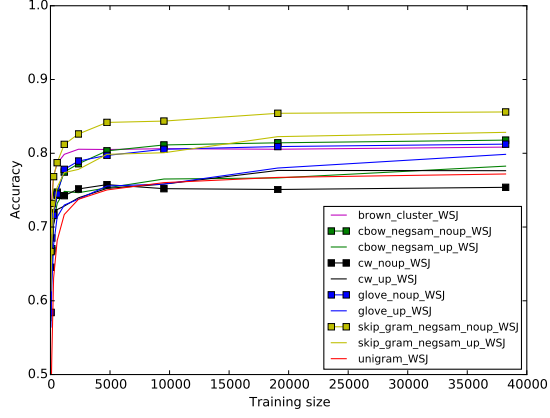
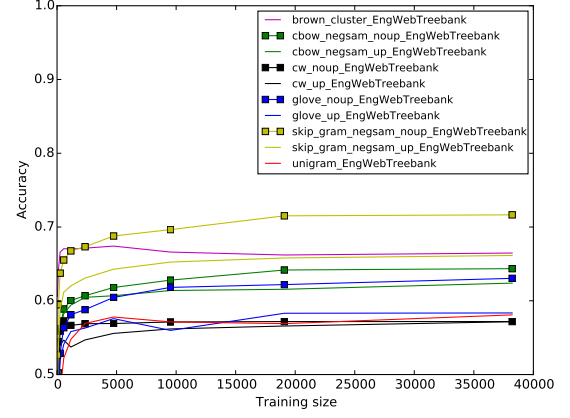


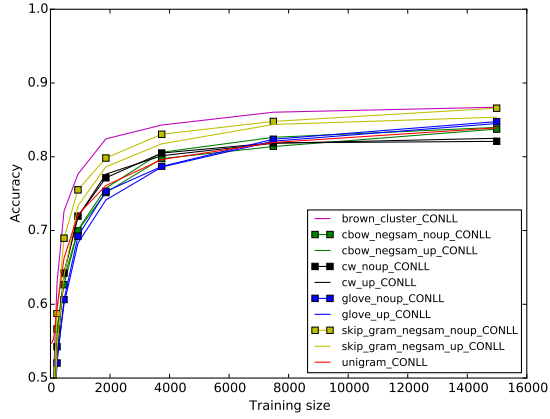
Figure 3: Out-of-vocabulary-words (OOV) accuracy for *in-domain* and *out-of-domain* test sets



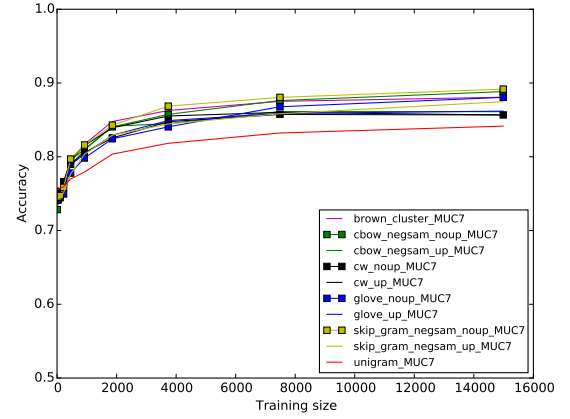
(a) POS-Tagging accuracy for *in domain* OOV



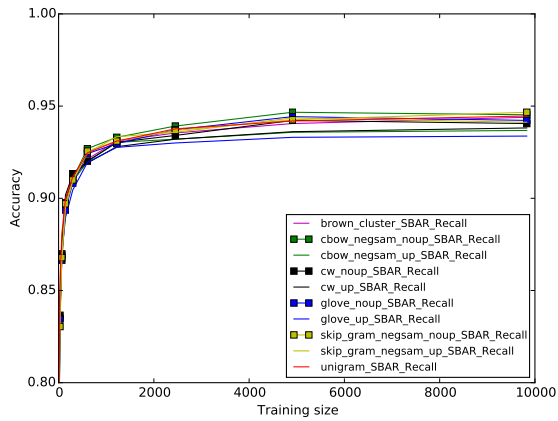
(b) POS-Tagging accuracy for *out domain* OOV



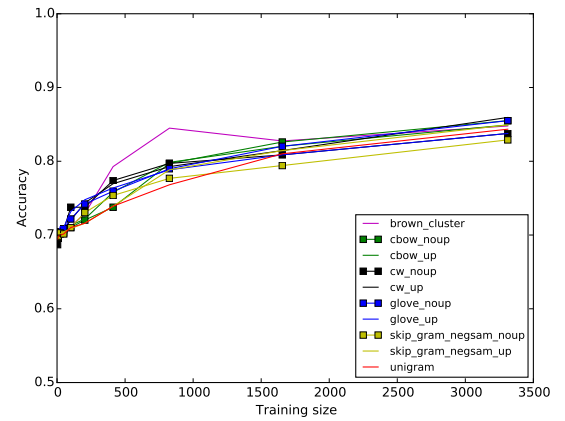
(c) NER accuracy for *out domain* OOV



(d) NER accuracy for *out domain* OOV



(e) Chunking accuracy for OOV *in domain* OOV



(f) MWE accuracy for *in domain* OOV

but the difference is marginal. Nevertheless, we found that fine-tuning can result in over-fitting, when the ones learned unsupervisedly were already good. Finally, by using word embeddings as features for MWE identification, we outperformed the state-of-the-art system. We could not find any trend that suggest that a word embedding method is better than other, thus suggesting that ... Future studies include learning representation of complex sequences such as MWEs.

Acknowledgments

Anonymised
Anonymised
Anonymised
Anonymised
Anonymised
Anonymised

References

- Rie Kubota Ando and Tong Zhang. 2005. A framework for learning predictive structures from multiple tasks and unlabeled data. *J. Mach. Learn. Res.*, 6:1817–1853, December.
- Jacob Andreas and Dan Klein. 2014. How much do word embeddings encode about syntax? In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 822–827, Baltimore, USA.
- Mohit Bansal, Kevin Gimpel, and Karen Livescu. 2014. Tailoring continuous word representations for dependency parsing. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 809–815, Baltimore, USA.
- Marco Baroni, Georgiana Dinu, and Germán Kruszewski. 2014. Don’t count, predict! a systematic comparison of context-counting vs. context-predicting semantic vectors. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 238–247, Baltimore, Maryland, June. Association for Computational Linguistics.
- James Bergstra and Yoshua Bengio. 2012. Random search for hyper-parameter optimization. *The Journal of Machine Learning Research*, 13(1):281–305.
- Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. 2011. Natural language processing (almost) from scratch. *The Journal of Machine Learning Research*, 12:2493–2537.
- John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *The Journal of Machine Learning Research*, 12:2121–2159.
- Guillaume Obozinski Edouard Grave and Francis Bach. 2013. Hidden markov tree models for semantic class induction. In *Proceedings of CoNLL*.
- J.R. Firth. 1957. A synopsis of linguistic theory 1930–1955. *Studies in linguistic analysis*, pages 1–32.
- Gholamreza Haffari, Marzieh Razavi, and Anoop Sarkar. 2011. An ensemble model that combines syntactic and semantic clustering for discriminative dependency parsing. In *ACL (Short Papers)*, pages 710–714. The Association for Computer Linguistics.
- Terry Koo, Xavier Carreras, and Michael Collins. 2008. Simple semi-supervised dependency parsing. In *In Proc. ACL/HLT*.
- Terry Koo, Alexander M. Rush, Michael Collins, Tommi Jaakkola, and David Sontag. 2010. Dual decomposition for parsing with non-projective head automata. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing, EMNLP ’10*, pages 1288–1298, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Olutobi Owoputi, Brendan O’Connor, Chris Dyer, Kevin Gimpel, Nathan Schneider, and Noah A Smith. 2013. Improved part-of-speech tagging for online conversational text with word clusters. In *HLT-NAACL*, pages 380–390.
- Nathan Schneider, Emily Danchik, Chris Dyer, and Noah A. Smith. 2014a. Discriminative lexical semantic segmentation with gaps: Running the mwe gamut. *Transactions of the Association of Computational Linguistics*, 2(1):193–206.
- Nathan Schneider, Spencer Onuffer, Nora Kazour, Emily Danchik, Michael T. Mordowanec, Henrietta Conrad, and Noah A. Smith. 2014b. Comprehensive annotation of multiword expressions in a social web corpus. In Nicoletta Calzolari, Khalid Choukri, Thierry Declerck, Hrafn Loftsson, Bente Maegaard, Joseph Mariani, Asuncion Moreno, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the Ninth International Conference on Language Resources and Evaluation*, pages 455–461, Reykjavík, Iceland, May. ELRA.
- Fei Sha and Fernando Pereira. 2003. Shallow parsing with conditional random fields. In *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology - Volume 1, NAACL ’03*, pages 134–141, Stroudsburg, PA, USA. Association for Computational Linguistics.

Valentin I. Spitkovsky, Hiyan Alshawi, Angel X. Chang, and Daniel Jurafsky. 2011. Unsupervised dependency parsing without gold part-of-speech tags. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing, EMNLP '11*, pages 1281–1290, Stroudsburg, PA, USA. Association for Computational Linguistics.

Kristina Toutanova, Dan Klein, Christopher D. Manning, and Yoram Singer. 2003. Feature-rich part-of-speech tagging with a cyclic dependency network. In *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology - Volume 1, NAACL '03*, pages 173–180, Stroudsburg, PA, USA. Association for Computational Linguistics.

Stephen Tratz and Eduard Hovy. 2011. A fast, accurate, non-projective, semantically-enriched parser. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing, EMNLP '11*, pages 1257–1268, Stroudsburg, PA, USA. Association for Computational Linguistics.

Joseph Turian, Lev Ratinov, and Yoshua Bengio. 2010. Word representations: a simple and general method for semi-supervised learning. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 384–394. Association for Computational Linguistics.