



## BIROn - Birkbeck Institutional Research Online

Singh, Manni and Weston, David and Levene, Mark (2020) Supervised phrase-boundary embeddings. In: Berthold, M.R. and Feelders, A. and Krempel, G. (eds.) Advances in Intelligent Data Analysis XVIII. Lecture Notes in Computer Science 12080. Springer, pp. 470-482. ISBN 9783030445836.

Downloaded from: <https://eprints.bbk.ac.uk/id/eprint/31525/>

*Usage Guidelines:*

Please refer to usage guidelines at <https://eprints.bbk.ac.uk/policies.html>  
contact [lib-eprints@bbk.ac.uk](mailto:lib-eprints@bbk.ac.uk).

or alternatively



# Supervised Phrase-Boundary Embeddings

Manni Singh<sup>(✉)</sup>, David Weston, and Mark Levene

Department of Computer Science and Information Systems,  
Birkbeck, University of London, London WC1E 7HX, UK  
{manni,dweston,mark}@dcs.bb.k.ac.uk

**Abstract.** We propose a new word embedding model, called SPhrase, that incorporates supervised phrase information. Our method modifies traditional word embeddings by ensuring that all target words in a phrase have exactly the same context. We demonstrate that including this information within a context window produces superior embeddings for both intrinsic evaluation tasks and downstream extrinsic tasks.

**Keywords:** Phrase embeddings · Named entity recognition · Natural language processing

## 1 Introduction

Word embeddings represent words with multidimensional vectors that are used in various models for applications such as, named entity recognition [9], query expansion [13], and sentiment analysis [21]. These embeddings are usually generated from a huge corpus with unsupervised learning models [3, 16, 18, 23, 24]. These models are based on describing target words by their neighbouring words which are also considered as contexts. The selection of these context words is generally linear (i.e.  $n$  words surrounding the target). Alternatively, arbitrary context words were used in [16] where context selection is based on the syntactic dependencies to the target word.

These models treat words as lexical units and create a context window surrounding a target word. This approach can be problematic when the context window for a target word contains only part of a phrase. For example, consider a scenario where a target word is close to (and to the right of) the named entity “George W. Bush” but the context window only retains the word “George”. Clearly this will generate ambiguity as the independent word “George” may refer another person (George Washington), location (George Street, Oxford) or a music band (George). To deal with the issue described above, [19] used a data-driven approach to identify and treat these phrases as individual tokens. While this technique may learn a phrase representation it cannot learn a representation of the individual words that comprise the phrase.

In our approach we obtain phrase information directly from Wikipedia. Terms from Wikipedia articles are formatted as hyperlinks to relevant articles. In a related method [22] these terms are extracted as named entities. This paper

interprets these terms as phrases. By using Wikipedia for phrase information (unlike [16]) we avoid needing additional grammatical information. This also gives us the potential to generate multi-lingual embeddings, although we do not pursue this here.

In this work, we are using phrase boundary information to generate word embedding in a non-compositional manner rather than a phrase embedding. We consider each of the words in the phrase as a part of the unit, where a unit can either be single word (i.e. not a link in the Wikipedia) or otherwise a bag of words. The embeddings are then learned for each of the unit members by considering surrounding units in the context.

In the following section we present related work in this domain, Sect. 3 presents our model and in Sects. 4 to 6 we give details of the implementation and the experiments.

## 2 Related Work

Word representations can be obtained from a language model where the goal is to predict a future word based on some previously observed information such as, a sentence, a sequence, or a phrase. For this task, various models can be utilised including: joint probabilities of observation that may include the Markov assumption. Under this assumption, we may say that the immediate future is independent of the entire past given the present. N-gram language models [4] use this assumption to predict token(s) using the previous  $N - 1$  tokens [17]. This can be constructed efficiently for very large datasets using neural network based language modelling (NNLM) [2].

The NNLM of [2] used a non-linear hidden layer between the input and output layers. A simpler network named the log bi-linear model was introduced in [20] by dropping the hidden layer between input and output layer. Instead of the hidden layer, context vectors were summed and projected to the output layer. This model was later used by [18] and named CBOW (Continuous Bag-of-words model), with a symmetric context (i.e. context words on both sides of the target word).

In addition, the Skip-gram model, was introduced in this work by reversing CBOW to predict context from the target word. Given a context range  $c$  and target word  $w_t$  the objective is to maximise the average log probability,

$$\sum_{-c \leq j \leq c} \log p(w_{t+j} | w_t)$$

The model defines  $p(w_{t+j} | w_t)$  using the softmax function,

$$p(w_O | w_I) = \frac{\exp(v'_{w_O} \top v_{w_I})}{\sum_{w=1}^W \exp(v'_w \top v_{w_I})}$$

where  $v_w$  and  $v'_w$  are the “input” and “output” vector representations of  $w$ , and  $W$  is the number of words in the vocabulary. However, due to the large vocabulary, the computation becomes impractical. Thus, Noise Contrastive Estimation (NCE) [7] was used that performs the same operation by sampling a very small amount of words  $k$  from the vocabulary as noise.

A similar technique is called Candidate Sampling [10] that combines noise samples with the true class, denoted as the set  $\mathcal{S}$ , with the objective to predict the true class from it, where  $Y$  is a set of true classes. Embeddings are scored as,

$$\hat{Y}_s = (X_s * W_s + b_s) - \log(E(s)).$$

Where  $X_s$  is a vector (embedding) corresponding to a word  $s \in \mathcal{S}$ ,  $W_s$  is the corresponding weight,  $b_s$  is the bias, and  $\mathbb{E}(s)$  is the expectation for  $s$ . Each score is approximated to a probability using the softmax function,

$$\text{Softmax}(\hat{Y}_s) = \frac{\exp \hat{Y}_s}{\sum_{s' \in \mathcal{S}} \exp \hat{Y}_{s'}}.$$

In addition to words, phrases may also be considered. In [18], the words comprising a phrase were joined using the delimiter ‘\_’ between them, and their joint embedding was learned. This scheme is called non-compositional embedding [8, 26]. Alternatively, compositional embeddings [8] are generated by merging word embeddings of phrase components using a composition function. The main difference in these schemes is that the previous learns the phrase embeddings while the latter just merges already learned word embeddings to make the phrase embeddings. Similarly, [3] introduced an extension of the Skip-gram model [18] that composes sub-word embeddings to make word embeddings with summation as the composition function.

### 3 The SPphrase Model

The proposed model uses information about which words belong to which phrases. This information can be conveniently represented as simply the locations for where phrases start and end, hence the name, *Supervised Phrase Boundary Representations model* (SPphrase).

The key assumption is that each word that comprises a phrase has the same context. This will produce an embedding where words that occur in the same phrase are likely to be close in the vector space. For example consider the sentence:

*British Airways to New York has Departed*

This sentence includes the (noun) phrase ‘New York’. Following the procedure for Word2vec we focus on the target word ‘New’ using a context window of size 1. The target, context pairs are (New, to) and (New, York). Repeating this procedure for the target word ‘York’, yields the target, context pairs (York, New) and (York, has).

For SPhrase, the context differs from Word2vec, both target words in ‘New York’ will have the same context based on the words immediately surrounding the phrase, hence the SPhrase target context pairs are (New, to), (New, has), (York, to), (York, has). Figure 1 highlights the context words for the word ‘New’ for both Word2vec and SPhrase.

Word2vec
British airways <b>to</b> <i>New</i> <b>York</b> has departed
SPhrase
British airways <b>to</b> <i>New</i> York <b>has</b> departed

**Fig. 1.** Context words for the target word *New* using Word2vec and SPhrase. The context words are in bold. The context size is 1.

In the above, we demonstrated the target context pairs induced by a target word that is a member of a phrase, where its context are individual words. In the following, we generalise the approach to handle the situation where phrases are part of a context. We do this by introducing the concept of a *unit*, where a unit consist of a sequence of words. A unit of length 1 represents individual words, a unit of length 2 represents two word phrases and so on for larger phrases.

Thus we measure the context simply in terms of units. Figure 2 provides an example of a context of size 2 each side. Note that the left context for SPhrase contains 3 words. Thus the context size measured in words will be larger for SPhrase than Word2vec if there is a phrase within the context window.

Word2vec
British <b>airways to</b> <i>Rome</i> <b>has departed</b>
SPhrase
<b>British airways to</b> <i>Rome</i> <b>has departed</b>

**Fig. 2.** Context words for the target word *Rome* using Word2vec and SPhrase. The context words are in bold. The context size is 2.

### 3.1 SPhrase Context Sampling

A standard approach to reduce the computation involved in generating embeddings is to shorten the effective context length by using only a sample of words from a context [18]. For SPhrase this can be achieved in several ways. First it can be done at the level of units not words, this is denoted *unit context sampling* (SPhrase). Second *random word context sampling* (R)<sup>1</sup> involves first performing unit context sampling, then for each unit that has a length greater than one only one word is sampled uniformly at random. This yields an effective context length that matches the context length of Word2vec. In addition to that, we generate embeddings named *without unit context sampling* (NU) where the target still is a unit but the context comprises individual words.

## 4 Methods and Datasets

### 4.1 Dataset

In order to generate an embedding using our approach, we require a corpus that has phrases annotated. Unfortunately this is not readily available, so we use a proxy for phrase annotation. In datasets that include hyperlinks we assume that the *hyperlink displayed text* is a phrase. One such data set is Wikipedia; we use the English Wikipedia dump version 20180920 that contains over 3 billion tokens. The proportion of tokens in phrases of length 2 is 2.5%; of length 3, 4, 5, and greater is respectively 0.8%, 0.3%, 0.2%, and less than 0.1%. Obviously not all phrases are represented as hyperlink text and not all hyperlink texts are phrases. Indeed the longest hyperlink text in our data set is of length 16,382 (it included internal formatting of Wikipedia). For our study we restricted maximum length to 10. The embedding vocabulary contained tokens with a frequency of at least 100 which gave us a total of 400,919 distinct tokens.

### 4.2 Parameter Settings

Training is performed in mini-batches of 60,000 tokens per batch with candidate sampling of 5000 classes per batch (value dictated by the available computational resource). The remaining parameters use standard values, the learning rate is initialised to 0.001 and optimisation is based on *Adam* optimiser [12] for stochastic learning. The learning decay is set to 10% (i.e. learning rate \* 0.9) after each epoch. The total number of the epochs is set to 20. The weighting scheme for selecting words in the context sampling is the same as for Word2vec [18].

## 5 Evaluation

There are two types of evaluation tasks commonly accepted: intrinsic and extrinsic. Intrinsic evaluation tasks determine the quality of embeddings. Under this

<sup>1</sup> Pretrained embeddings are available at: <https://github.com/ManniSingh/SPhrase>.

class, word similarity/relatedness tasks are generally based on cosine distance as a metric to find similarity between two word vectors. Extrinsic evaluation tasks, on the other hand, are based on specific downstream tasks such as, named entity recognition (NER), sentiment classification, topic detection. In this work, we are doing similarity based intrinsic evaluation and NER based extrinsic evaluation.

## 6 Experimental Design

### 6.1 Intrinsic Evaluation

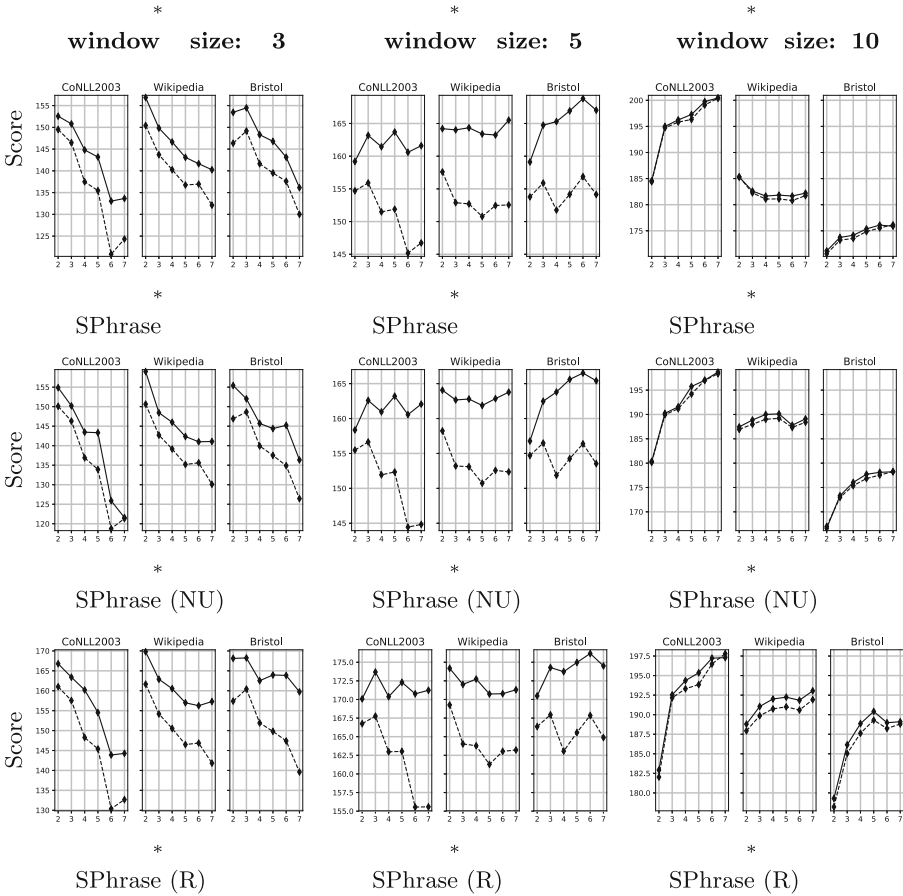
The following experiments fit into the so-called *intrinsic* category of embedding evaluation. We aim to demonstrate that although the total number of phrases in our dataset is small compared to the number of words, they do have a positive impact on the resulting embeddings. In order to determine an optimal configuration of the method, intrinsic evaluation is done on embeddings trained on the first 10% of the corpus; see Fig. 3. As a result, the extrinsic evaluation described Sect. 6.2, the performance of the optimal configuration in this evaluations is: SPhrase (R) with window size 5. For the extrinsic evaluation only the optimal configuration is used and the embeddings are trained on the full corpus.

In the following experiments we compare SPhrase embeddings with the ones generated by Word2vec. It is known that increasing the context window size generally improves the quality of the embedding. Recall that the expected context size for each target word is the same for Word2vec and SPhrase due to word context sampling.

We expect that words in phrases should be mapped to similar locations in the embedding, i.e. words within a phrase should be closer together than words that are not in the same phrase. In the following we first perform experiment on pairwise similarity and then we investigate further structure with an analogy task.

**Pairwise Similarity.** For pairwise similarity experiments we use phrases from three datasets.

- CoNLL-2003 English dataset [25]. From this dataset multi-word named entities were extracted. These are used as phrases, in total there are 12,999. The maximum phrase length is 7 in this dataset, so we restricted the following two datasets to this as well.
- From our Wikipedia training corpus we obtained 16,470 phrases from the first 1,000,000 tokens. This dataset comes from our training data, so we assume we should obtain good results in this case.
- Bristol [15] - from this dataset we selectively used the entity list and found 87,209 phrases.



**Fig. 3.** Similarity scores comparison for the phrases relative to 100 random words representing: *unit context sampling* (SPhrase), Without *unit context sampling* (NU) and, with *random word context sampling* (R). Where SPhrase (in bold) and Word2vec (dashed) are compared on phrase lengths 2–7 (in horizontal axis) with higher the score the better it performed.

In order to investigate how the distances of words within a phrase compare to distances of words with random words in the datasets we use the following,

$$\text{Similarity Score} = \frac{1}{N_i(l-1)} \sum_{i=1}^{l-1} b(w_i, w_{i+1}, r)$$

where,

$$b(w_i, w_{i+1}, r) = \begin{cases} 1 & s(w_i, w_{i+1}) > s(w_i, r), \\ 0 & \text{otherwise,} \end{cases}$$



where  $r$  is a word selected at random from another phrase. A new word is drawn for each phrase pair comparison. The similarity score is calculated 100 times and the overall average is taken in order to reduce the noise generated by selecting only one word for each comparison. The interpretation of this is similar to the cosine score in that the larger the value the better.

We computed scores for phrase lengths up to and including length 7. We have used context window sizes 3, 5 and 10. Figure 3 shows these scores for the context sampling regimes: with *unit context sampling*, without *unit context sampling*, and *word context sampling*.

We can see that regardless of the embedding, the scores in general reduce as the phrase gets longer. However, the larger the window size the more Word2vec and SPhrase agree. This is what we should expect, since there will be greater overlap in the context words between SPhrase and Word2vec. Nevertheless we see that, overall, SPhrase performs better.

**Google Analogy Test Set.** Analogy based tasks are widely used, e.g. [5, 6, 11] to evaluate the quality of word embeddings. One well known test set is the Google analogy test set [18]. This dataset comprises rows of four words, such as **known unknown informed uninformed**. The analogy task is to predict the final word using the first three using simple vector addition/subtraction of their vector representations. Informally the task attempts to show how well words follow the vector relationship

$$unknown - known = uninformed - informed$$

**Table 1.** Scores on Google analogy dataset with *unit context sampling* (SPhrase), here accuracy is the total correct count on the total count of instances.

	Accuracy - displayed to 3 decimal places						Count
	Window size 3		Window size 5		Window size 10		
	SPhrase	Word2vec	SPhrase	Word2vec	SPhrase	Word2vec	
<i>capital-world</i>	0.727	0.628	0.746	0.658	0.815	0.782	4524
<i>capital-common-countries</i>	0.872	0.848	0.941	0.856	0.976	0.941	506
<i>city-in-state</i>	0.660	0.480	0.715	0.583	0.647	0.677	2467
gram3-comparative	0.848	0.806	0.758	0.813	0.643	0.670	1332
gram2-opposite	0.223	0.220	0.220	0.222	0.206	0.204	812
gram8-plural	0.755	0.736	0.715	0.744	0.641	0.727	1332
gram4-superlative	0.379	0.396	0.345	0.366	0.279	0.262	1122
gram9-plural-verbs	0.639	0.559	0.536	0.546	0.453	0.521	870
gram6-nationality-adjective	0.846	0.784	0.838	0.815	0.854	0.853	1599
family	0.603	0.595	0.595	0.638	0.581	0.543	506
gram7-past-tense	0.472	0.515	0.474	0.492	0.441	0.470	1560
currency	0.047	0.042	0.021	0.021	0.018	0.016	866
gram1-adjective-to-adverb	0.104	0.087	0.119	0.121	0.132	0.148	992
gram5-present-participle	0.517	0.520	0.509	0.486	0.479	0.455	1056
all	0.601	0.545	0.597	0.565	0.581	0.587	19544

**Table 2.** Scores on Google analogy dataset without *unit context sampling* (NU), here accuracy is the total correct count on the total count of instances.

	Accuracy - displayed to 3 decimal places						Count
	Window size 3		Window size 5		Window size 10		
	SPhrase	Word2vec	SPhrase	Word2vec	SPhrase	Word2vec	
<i>capital-world</i>	0.671	0.628	0.725	0.658	0.744	0.782	4524
<i>capital-common-countries</i>	0.881	0.848	0.935	0.856	0.929	0.941	506
<i>city-in-state</i>	0.653	0.480	0.645	0.583	0.652	0.677	2467
gram3-comparative	0.706	0.806	0.696	0.813	0.519	0.670	1332
gram2-opposite	0.217	0.220	0.197	0.222	0.172	0.204	812
gram8-plural	0.726	0.736	0.712	0.744	0.661	0.727	1332
gram4-superlative	0.273	0.396	0.298	0.366	0.269	0.262	1122
gram9-plural-verbs	0.577	0.559	0.548	0.546	0.477	0.521	870
gram6-nationality-adjective	0.855	0.784	0.821	0.815	0.827	0.853	1599
family	0.569	0.595	0.553	0.638	0.502	0.543	506
gram7-past-tense	0.453	0.515	0.483	0.492	0.414	0.470	1560
currency	0.039	0.042	0.024	0.021	0.028	0.016	866
gram1-adjective-to-adverb	0.130	0.087	0.173	0.121	0.168	0.148	992
gram5-present-participle	0.511	0.520	0.509	0.486	0.492	0.455	1056
all	0.565	0.545	0.576	0.565	0.553	0.587	19544

The dataset is divided into categories, some of which are inherently phrase-based. In the category `capital-common-countries` a typical line is:

`Athens Greece Baghdad Iraq`

Both *Athens Greece* and *Baghdad Iraq* can be reasonably construed to be phrases, unlike in the first example above. Two other categories have this same character, namely `capital-world` and `city-in-state`.

Example rows are: `Athens Greece Canberra Australia` and `Chicago Illinois Houston Texas` respectively.

With this in mind we show the accuracy of SPhrase and Word2vec stratified by category, in addition to the overall accuracy that is usually reported. The categories that have a phrasal quality are italicised in Tables 1, 2 and 3. We see that, overall, SPhrase performs better in these categories.

## 6.2 Extrinsic Evaluation

We use Conll2003 English [25] and Wikigold [1] to evaluate the performance of the embeddings generated. The Conll dataset is widely used to evaluate various NER based models. It contains 203,621 tokens in the training set, while validation and test set contains 51,362 and 46,435 tokens respectively. On the other hand, Wikigold provides a single data file of 39,007 tokens that we used for testing while the NER models were trained with Conll train and validation data. We used SPhrase (R) model with window size 5 since this configuration demonstrated significant improvements over Word2vec as shown in Fig. 3. We recreated the BLSTMs and CRF based model [14] but without any feature engineering.

**Table 3.** Scores on Google analogy dataset with *random word context sampling* (R), here accuracy is the total correct count on the total count of instances.

	Accuracy - displayed to 3 decimal places						Count
	Window size 3		Window size 5		Window size 10		
	SPhrase	Word2vec	SPhrase	Word2vec	SPhrase	Word2vec	
<i>capital-world</i>	0.637	0.628	0.718	0.658	0.766	0.782	4524
<i>capital-common-countries</i>	0.858	0.848	0.903	0.856	0.953	0.941	506
<i>city-in-state</i>	0.664	0.480	0.623	0.583	0.663	0.677	2467
gram3-comparative	0.845	0.806	0.803	0.813	0.682	0.670	1332
gram2-opposite	0.224	0.220	0.245	0.222	0.196	0.204	812
gram8-plural	0.772	0.736	0.731	0.744	0.655	0.727	1332
gram4-superlative	0.373	0.396	0.392	0.366	0.257	0.262	1122
gram9-plural-verbs	0.575	0.559	0.586	0.546	0.474	0.521	870
gram6-nationality-adjective	0.818	0.784	0.824	0.815	0.831	0.853	1599
family	0.615	0.595	0.581	0.638	0.595	0.543	506
gram7-past-tense	0.479	0.515	0.520	0.492	0.460	0.470	1560
currency	0.040	0.042	0.024	0.021	0.023	0.016	866
gram1-adjective-to-adverb	0.090	0.087	0.127	0.121	0.172	0.148	992
gram5-present-participle	0.526	0.520	0.455	0.486	0.479	0.455	1056
all	0.576	0.545	0.588	0.565	0.576	0.587	19544

**Table 4.** Comparison of Word2vec with SPhrase(NU) on Conll2003 English and Wikigold dataset

Model	Conll2003Eng	Wikigold
Word2Vec	83.82 $\pm$ 0.3831	55.49 $\pm$ 0.4708
SPhrase	<b>88.93 <math>\pm</math> 0.1115</b>	<b>66.01 <math>\pm</math> 0.4172</b>

We trained this in 20 epochs with evaluating on validation data each time. We performed 10 instances for each of these models and presented the range of F1 scores (using Conll2003 evaluation script). Table 4 displays the results that show a significant improvement over the Word2vec model trained on the same corpus.

## 7 Concluding Remarks

This investigation demonstrates that using phrasal information can directly enrich word embeddings. In this work, we presented an alternative context sampling technique to that used in skip-gram Word2vec. We note that the SPhrase approach is not limited to augmenting Word2Vec, it can also be applied to morphological extensions such as Fasttext [3].

We used the displayed text from hyperlinks as a proxy for phrases, and in this sense SPhrase is supervised. We are, however, planning to generalise the methodology by investigating whether we can identify useful phrase boundaries in a completely unsupervised fashion.

## References

1. Balasuriya, D., Ringland, N., Nothman, J., Murphy, T., Curran, J.R.: Named entity recognition in Wikipedia. In: Proceedings of the 2009 Workshop on the People's Web Meets NLP: Collaboratively Constructed Semantic Resources. People's Web 2009, pp. 10–18. Association for Computational Linguistics, Stroudsburg, PA, USA (2009). <http://dl.acm.org/citation.cfm?id=1699765.1699767>
2. Bengio, Y., Ducharme, R., Vincent, P., Jauvin, C.: A neural probabilistic language model. *J. Mach. Learn. Res.* **3**(Feb), 1137–1155 (2003)
3. Bojanowski, P., Grave, E., Joulin, A., Mikolov, T.: Enriching word vectors with subword information. *Trans. Assoc. Comput. Linguist. (TACL)* **5**(1), 135–146 (2017). <http://www.aclweb.org/anthology/Q17-1010>
4. Brants, T., Papat, A.C., Xu, P., Och, F.J., Dean, J.: Large language models in machine translation. In: Proceedings of the Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning, pp. 858–867 (2007)
5. Bruni, E., Tran, N.K., Baroni, M.: Multimodal distributional semantics. *J. Arti. Intell. Res.* **49**(1), 1–47 (2014). <http://dl.acm.org/citation.cfm?id=2655713.2655714>
6. Finkelstein, L., et al.: Placing search in context: the concept revisited. In: Proceedings of the 10th International Conference on World Wide Web, WWW 2001, pp. 406–414. ACM, New York (2001). <https://doi.org/10.1145/371920.372094>
7. Gutmann, M.U., Hyvärinen, A.: Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics. *J. Mach. Learn. Res.* **13**(Feb), 307–361 (2012)
8. Hashimoto, K., Tsuruoka, Y.: Adaptive joint learning of compositional and non-compositional phrase embeddings. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 205–215. Association for Computational Linguistics, Berlin, Germany, August 2016 (2016). <https://doi.org/10.18653/v1/P16-1020>, <http://www.aclweb.org/anthology/P16-1020>
9. Huang, Z., Xu, W., Yu, K.: Bidirectional LSTM-CRF models for sequence tagging. arXiv preprint [arXiv:1508.01991](https://arxiv.org/abs/1508.01991) abs/1508.01991 (2015)
10. Jean, S., Cho, K., Memisevic, R., Bengio, Y.: On using very large target vocabulary for neural machine translation. In: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pp. 1–10. Association for Computational Linguistics (2015)
11. Jurgens, D.A., Turney, P.D., Mohammad, S.M., Holyoak, K.J.: Semeval-2012 task 2: measuring degrees of relational similarity. In: Proceedings of the First Joint Conference on Lexical and Computational Semantics - Volume 1: Proceedings of the Main Conference and the Shared Task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation, SemEval 2012, pp. 356–364. Association for Computational Linguistics, Stroudsburg, PA, USA (2012), <http://dl.acm.org/citation.cfm?id=2387636.2387693>
12. Kingma, D.P., Ba, J.: Adam: a method for stochastic optimization. arXiv abs/1412.6980 (2014)
13. Kuzi, S., Shtok, A., Kurland, O.: Query expansion using word embeddings. In: Proceedings of the 25th ACM International on Conference on Information and Knowledge Management, CIKM 2016, pp. 1929–1932. ACM, New York (2016)

14. Lample, G., Ballesteros, M., Subramanian, S., Kawakami, K., Dyer, C.: Neural architectures for named entity recognition. In: NAACL HLT 2016, the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, 12–17 June 2016, pp. 260–270 (2016)
15. Lansdall-Welfare, T., Sudhahar, S., Thompson, J., Lewis, J., Team, F.N., Cristianini, N.: Content analysis of 150 years of british periodicals. *Proc. Nat. Acad. Sci.* **114**(4), E457–E465 (2017)
16. Levy, O., Goldberg, Y.: Dependency-based word embeddings. In: Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), vol. 2, pp. 302–308 (2014)
17. Martin, J.H., Jurafsky, D.: *Speech and Language Processing: an Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. Pearson/Prentice Hall, Upper Saddle River (2009)
18. Mikolov, T., Chen, K., Corrado, G., Dean, J.: Efficient estimation of word representations in vector space. arXiv abs/1301.3781 (2013). <http://arxiv.org/abs/1301.3781>
19. Mikolov, T., Sutskever, I., Chen, K., Corrado, G., Dean, J.: Distributed representations of words and phrases and their compositionality. In: Proceedings of the 26th International Conference on Neural Information Processing Systems (NIPS), NIPS 2013, vol. 2, pp. 3111–3119. Curran Associates Inc., USA (2013)
20. Mnih, A., Hinton, G.: Three new graphical models for statistical language modelling. In: Proceedings of the 24th International Conference on Machine Learning. ICML 2007, pp. 641–648. ACM, New York (2007)
21. Nakov, P., Ritter, A., Rosenthal, S., Sebastiani, F., Stoyanov, V.: SemEval-2016 task 4: sentiment analysis in Twitter. In: Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016), pp. 1–18, USA (2016)
22. Nothman, J., Curran, J.R., Murphy, T.: Transforming Wikipedia into named entity training data. In: Proceedings of the Australasian Language Technology Association Workshop 2008, pp. 124–132, Hobart, Australia, December 2008. <http://www.aclweb.org/anthology/U08-1016>
23. Pennington, J., Socher, R., Manning, C.D.: Glove: global vectors for word representation. In: Empirical Methods in Natural Language Processing (EMNLP), pp. 1532–1543 (2014). <http://www.aclweb.org/anthology/D14-1162>
24. Salle, A., Villavicencio, A., Idiart, M.: Matrix factorization using window sampling and negative sampling for improved word representations. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, Berlin, Germany, 7–12 August 2016, vol. 2, Short Papers (2016)
25. Tjong Kim Sang, E.F., De Meulder, F.: Introduction to the CoNLL-2003 shared task: language-independent named entity recognition. In: Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003, vol. 4, pp. 142–147. Association for Computational Linguistics, Japan (2003)
26. Yu, M., Dredze, M.: Learning composition models for phrase embeddings. *Trans. Assoc. Comput. Linguist. (TACL)* **3**(1), 227–242 (2015). <http://www.aclweb.org/anthology/Q15-1017>

**Open Access** This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

