

Big Data Small Data, In Domain Out-of Domain, Known Word Unknown Word: The Impact of Word Representation on Sequence Labelling Tasks

Lizhen Qu^{1,2}, Gabriela Ferraro^{1,2}, Liyuan Zhou¹, Weiwei Hou¹, Nathan Schneider³ and Timothy Baldwin⁴

¹ NICTA / Locked Bag 8001, ACT 2601, Australia

² The Australian National University

³ The University of Melbourne, VIC 3010, Australia

⁴ University of Edinburgh, EH8 9AB, UK.

{lizhen.qu, gabriela.ferraro, liyuan.zho, weiwei.hou}@nicta.com.au
nschneid@cs.cmu.edu
tb@ldwin.net

Abstract

Word embeddings – distributed word representations that can be learned from unlabelled data – have been shown to have high utility in many natural language processing applications. In this paper, we perform an extrinsic evaluation of five popular word embedding methods in the context of four sequence labelling tasks: POS-tagging, syntactic chunking, NER and MWE identification. A particular focus of the paper is analysing the effects of task-based updating of word representations. We show that when using word embeddings as features, as few as several hundred training instances are sufficient to achieve competitive results, and that word embeddings lead to improvements over OOV words and out of domain. Perhaps more surprisingly, our results indicate there is little difference between the different word embedding methods, and that simple Brown clusters are often competitive with word embeddings across all tasks we consider.

1 Introduction

Recently, distributed word representations have grown to become a mainstay of natural language processing (NLP), and been shown to have empirical utility in a myriad of tasks (e.g., ?; ?; ?; ?). The underlying idea behind distributed word representations is simple: to map each word w in our vocabulary V onto a continuous-valued vector of dimensionality $d \ll |V|$. Words that are similar (e.g., with respect to syntax or lexical semantics) will ideally be mapped to similar regions of the vector space, implicitly supporting both generalisation across in-vocabulary (IV) items, and counteracting the effects of data sparsity for low-frequency and out-of-vocabulary (OOV) items.

Without some means of automatically deriving the vector representations without reliance on labelled data, however, word embeddings would have little practical utility. Fortunately, it has been shown that they can be “pre-trained” from unlabelled text data using various algorithms to model the distributional hypothesis (i.e., that words which occur in similar contexts tend to be semantically similar). Pre-training methods have been refined considerably in recent years, and scaled up to increasingly large corpora.

As with other machine learning methods, it is well known that the quality of the pre-trained word embeddings depends heavily on factors including parameter optimisation, the size of the training data, and the fit with the target application. For example, (Mikolov et al., 2013) showed that the optimal dimensionality for word embeddings is task-specific. One factor which has received relatively little attention in NLP is the effect of “updating” the pre-trained word embeddings as part of the task-specific training, based on self-taught learning (Liu et al., 2012). Updating leads to word representations that are task-specific, but often at the cost of over-fitting low-frequency and OOV words.

In this paper, we perform an extensive evaluation of five word embedding approaches under fixed experimental conditions, applied to four sequence labelling tasks: POS-tagging, full-text chunking, named entity recognition (NER), and multiword expression (MWE) identification. In this, we explore the following research questions:

RQ1: are word embeddings better than baseline approaches of one-hot unigram features and Brown clusters?

RQ2: do word embeddings require less training data (i.e. generalise better) than one-hot unigram features?

RQ3: what is the impact of updating word embeddings in sequence labeling tasks, both empirically over the target task and geo-

metrically over the vectors?

RQ4: what is the impact of word embeddings (with and without updating) on both OOV items (relative to the training data) and out-of-domain data?

RQ5: overall, are some word embeddings better than others in a sequence labelling context?

2 Word Representations

2.1 Types of Word Representations

?) identifies three varieties of word representations: *distributional*, *cluster-based*, and *distributed*.

Distributional representation methods map each word w to a context word vector \mathbf{C}_w , which is constructed directly from co-occurrence counts between w and its context words. The learning methods either store the co-occurrence counts between two words w and i directly in C_{wi} (?: ?; ?) or project the concurrence counts between words into a lower dimensional space (?: ?), using dimensionality reduction techniques such as SVD (?) and LDA (?).

Cluster-based representation methods build clusters of words by applying either soft or hard clustering algorithms (?: ?). Some of them also rely on a co-occurrence matrix of words (?). The Brown clustering algorithm (?) is the best-known method in this category.

Distributed representation methods usually map words into dense, low-dimensional, continuous-valued vectors, with $\mathbf{x} \in R^d$, where d is referred to as the word dimension.

2.2 Selected Word Representations

Over a range of sequence labelling tasks, we evaluate five methods for inducing word representations: Brown clustering (?) (“BROWN”), the neural language model of Collobert & Weston (“CW”) (?), the continuous bag-of-words model (“CBOW”) (?), the continuous skip-gram model (“SKIP-GRAM”) (?), and Global vectors (“GLOVE”) (?). With the exception of CW, all have have been shown to be at or near state-of-the-art in recent empirical studies (?: ?). CW is included because it was highly influential in earlier research, and the pre-trained embeddings are still used to some degree in NLP. The training of these word representations is unsupervised: the common underlying idea is to predict occurrence of words in the neighbouring context. Their train-

ing objectives share the same form, which is a sum of local training factors $J(w, \text{ctx}(w))$,

$$L = \sum_{w \in V} J(w, \text{ctx}(w))$$

where V is the vocabulary of a given corpus, and $\text{ctx}(w)$ denotes the local context of word w . The local context of a word can either be its previous k words, or the k words surrounding it. Local training factors are designed to capture the relationship between w and its local contexts of use, either by predicting w based on its local context, or using w to predict the context words. Other than BROWN, which utilises a cluster-based representation, all the other methods employ a distributed representation.

The starting point for CBOW and SKIP-GRAM is to employ softmax to predict word occurrence:

$$J(w, \text{ctx}(w)) = -\log \left(\frac{\exp(\mathbf{v}_w^T \mathbf{v}_{\text{ctx}(w)})}{\sum_{j \in V} \exp(\mathbf{v}_j^T \mathbf{v}_{\text{ctx}(w)})} \right)$$

where $\mathbf{v}_{\text{ctx}(w)}$ denotes the distributed representation of the local context of word w . CBOW derives $\mathbf{v}_{\text{ctx}(w)}$ based on averaging over the context words. That is, it estimates the probability of each w given its local context. In contrast, SKIP-GRAM applies softmax to each context word of a given occurrence of word w . In this case, $\mathbf{v}_{\text{ctx}(w)}$ corresponds to the representation of one of its context words. This model can be characterised as predicting context words based on w . In practice, softmax is too expensive to compute over large corpora, and thus ?) use hierarchical softmax and negative sampling to scale up the training.

CW considers the local context of a word w to be m words to the left and m words to the right of w . The concatenation of the embeddings of w and all its context words are taken as input to a neural network with one hidden layer, which produces a higher level representation $f(w) \in R^d$. Then the learning procedure replaces the embedding of w with that of a randomly sampled word w' and generates a second representation $f(w') \in R^d$ with the same neural network. The training objective is to maximise the difference between them:

$$J(w, \text{ctx}(w)) = \max(0, 1 - f(w) + f(w'))$$

This approach can be regarded as negative sampling with only one negative example.

GLOVE assumes the dot product of two word embeddings should be similar to the logarithm of

Data set	Size	Words
UMBC	48.1GB	3G
One Billion	4.1GB	1G
English Wikipedia	49.6GB	3G

Table 1: Corpora used to pre-train the word embeddings

the co-occurrence count X_{ij} of the two words. As such, the local factor $J(w, \text{ctx}(w))$ becomes:

$$g(X_{ij})(\mathbf{v}_i^T \mathbf{v}_j + b_i + b_j - \log(X_{ij}))^2$$

where b_i and b_j are the bias terms of words i and j , respectively, and $g(X_{ij})$ is a weighting function based on the co-occurrence count. This weighting function controls the degree of agreement between the parametric function $\mathbf{v}_i^T \mathbf{v}_j + b_i + b_j$ and $\log(X_{ij})$. Frequently co-occurring word pairs will be larger weight than infrequent pairs, up to a threshold.

BROWN partitions words into a finite set of word classes V . The conditional probability of seeing the next word is defined to be:

$$p(w_k | w_{k-m}^{k-1}) = p(w_k | h_k) p(h_k | h_{k-m}^{k-1})$$

where h_k denotes the word class of the word w_k , w_{k-m}^{k-1} are the previous m words and h_{k-m}^{k-1} are their respective word classes. Then $J(w, \text{ctx}(w)) = -\log p(w_k | w_{k-m}^{k-1})$. Since there is no tractable method to find an optimal partition of word classes, the method uses only a bi-gram class model, and utilises hierarchical clustering as an approximation method to find a sufficiently good partition of words.

2.3 Building Word Representations

For a fair comparison, we train BROWN, CBOW, SKIP-GRAM, and GLOVE on a fixed corpus, comprised of freely available corpora, as detailed in Tab. 1. The joint corpus was preprocessed with the Stanford CoreNLP sentence splitter and tokeniser. 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 NUM2.NUM2). Due to the computational complexity of the pre-training, for CW, we simply downloaded the pre-compiled embeddings from: <http://metaoptimize.com/projects/wordreprs>.

The dimensionality of the word embeddings and the size of the context window are the key hyperparameters when learning distributed representations. We use all combinations of the following

values to train word embeddings on the combined corpus:

- **Embedding dim.** $d \in \{25, 50, 100, 200\}$
- **Context window size** $m \in \{1, 5, 10\}$

BROWN requires only the number of clusters as a hyperparameter. We perform clustering with $b \in \{250, 500, 1000, 2000, 4000\}$ clusters.

3 Sequence Labelling Tasks

We evaluate the different word representations over four sequence labelling tasks: POS-tagging (POS-tagging), full-text chunking (Chunking), NER (NER) and MWE identification (MWE). For each task, we fed features into a first order linear-chain graph transformer (?) made up of two layers: the upper layer is identical to a linear-chain CRF (?), and the lower layer consists of word representation and hand-crafted features. If we treat word representations as fixed, the graph transformer is a simple linear-chain CRF. On the other hand, if we can treat the word representations as model parameters, the model is equivalent to a neural network with word embeddings as the input layer. We trained all models using AdaGrad (?).

As in (?), at each word position, we construct word representation features from the words in a context window of size two to either side of the target word, based on the pre-trained representation of each word type. For BROWN, the features are the prefix features extracted from word clusters in the same way as (?). As a baseline (and to test RQ1), we include a one-hot representation (which is equivalent to a linear-chain CRF with only lexical context features).

Our hand-crafted features for POS-tagging, Chunking and MWE, are those used by (?), (?) and (?), respectively. For NER, we use the same feature space as (?), except for the previous two predictions, because we want to evaluate all word representations with the same type of model – a first-order graph transformer.

In training the distributed word representations, we consider two settings: (1) the word representations are fixed during sequence model training; and (2) the graph transformer updated the token-level word representations during training.

As outlined in Tab. 2, for each sequence labelling task, we experiment over the de facto corpus, based on pre-existing training–dev–test splits

	Training	Development	In-domain Test	Out-of-domain Test
POS-tagging	WSJ Sec. 0-18	WSJ Sec. 19-21	WSJ Sec. 22-24	EWT
Chunking	WSJ	WSJ (1K sentences)	WSJ (CoNLL-00 test)	Brown
NER	Reuters (CoNLL-03 train)	Reuters (CoNLL-03 dev)	Reuters (CoNLL-03 test)	MUC7
MWE	EWT (500 docs)	EWT (100 docs)	EWT (123 docs)	—

Table 2: Training, development and test (in- and out-of-domain) data for each sequence labelling task.

where available:¹

POS-tagging: the Wall Street Journal portion of the Penn Treebank (?): “WSJ”) with Penn POS tags

Chunking: the Wall Street Journal portion of the Penn Treebank (“WSJ”), converted into IOB-style full-text chunks using the CoNLL conversion scripts for training and dev, and the WSJ-derived CoNLL-2000 full text chunking test data for testing (?)

NER: the English portion of the CoNLL-2003 English Named Entity Recognition data set, for which the source data was taken from Reuters newswire articles (?): “Reuters”)

MWE: the MWE dataset of (?), over a portion of text from the English Web Treebank² (“EWT”)

For all tasks other than MWE,³ we additionally have an out-of-domain test set, in order to evaluate the out-of-domain robustness of the different word representations, with and without updating. These datasets are as follows:

POS-tagging: the English Web Treebank with Penn POS tags (“EWT”)

Chunking: the Brown Corpus portion of the Penn Treebank (“Brown”), converted into IOB-style full-text chunks using the CoNLL conversion scripts

NER: the MUC-7 named entity recognition corpus⁴ (“MUC7”)

For reproducibility, we tuned the hyperparameters with random search over the development data for each task (?). In this, we randomly sampled 50 distinct hyperparameter sets with the same random seed for the non-updating models (i.e. the models that don’t update the word representation), and sampled 100 distinct hyperparameter sets for the updating models (i.e. the models that do). For

each set of hyperparameters and task, we train a model over its training set and choose the best one based on its performance on development data (?). We also tune the word representation hyperparameters – namely, the word vector size d and context window size m (distributed representations), and in the case of Brown, the number of clusters.

For the updating models, we found that the results over the test data were always inferior to those that do not update the word representations, due to the higher number of hyperparameters and small sample size (i.e. 100). Since the two-layer model of the graph transformer contains a distinct set of hyperparameters for each layer, we reuse the best-performing hyperparameter settings from the non-updating models, and only tune the hyperparameters of AdaGrad for the word representation layer. This method requires only 32 additional runs and achieves consistently better results than 100 random draws.

In order to test the impact of the volume of training data on the different models (**RQ2**), we split the training set into 10 partitions based on a base-2 log scale (i.e., the second smallest partition will be twice the size of the smallest partition), and created 10 successively larger training sets by merging these partitions from the smallest one to the largest one, and used each of these to train a model. From these, we construct learning curves over each task.

For ease of comparison with previous results, we evaluate both in- and out-of-domain using chunk/entity/expression-level F1-measure (“F1”) for all tasks except POS-tagging, for which we use token-level accuracy (“ACC”). To test performance over OOV (unknown) tokens – i.e., the words that do not occur in the training set – we use token-level accuracy for all tasks (e.g. for Chunking, we evaluate whether the full IOB tag is correct or not), due to the sparsity of all-OOV chunks/NEs/MWEs.

4 Experimental Results and Discussion

We structure our evaluation by stepping through each of our five research questions (**RQ1–5**) from

¹For the MWE dataset, no such split pre-existed, so we constructed our own.

²<https://catalog.ldc.upenn.edu/LDC2012T13>

³Unfortunately, there is no second domain which has been hand-tagged with MWEs using the method of (?) to use as an out-of-domain test corpus.

⁴<https://catalog.ldc.upenn.edu/LDC2001T02>

Task	Benchmark	<i>In-domain</i> Test set	<i>Out-of-domain</i> Test set
POS-tagging (ACC)	0.972 (?)	0.959 (SKIP-GRAM+UP)	0.910 (SKIP-GRAM)
Chunking (F1)	0.942 (?)	0.938 (BROWN _{b=2000})	0.676 (GLOVE)
NER (F1)	0.893 (?)	0.868 (SKIP-GRAM)	0.736 (SKIP-GRAM)
MWE (F1)	0.625 (?)	0.654 (CW)	—

Table 3: State-of-the-art results vs. our best results for in-domain and out-of-domain test sets.

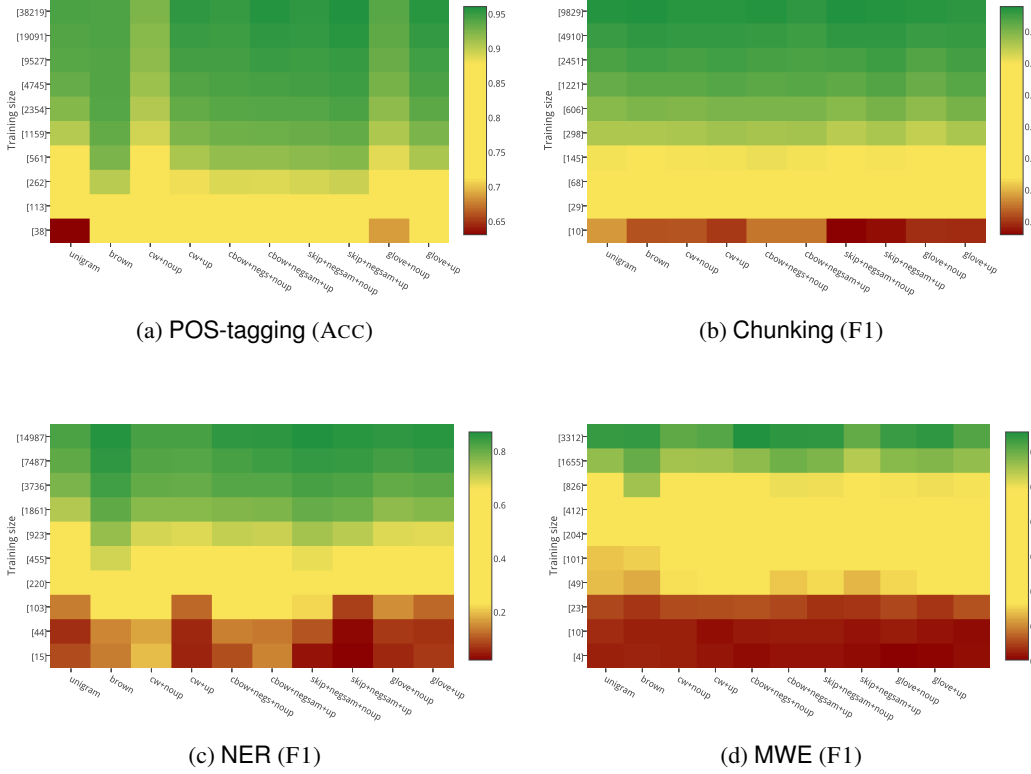


Figure 1: Results for each type of word representation over POS-tagging, Chunking, NER and MWE, optionally with updating (“+UP”). The y -axis indicates the training data sizes (on a log scale). Green = high performance, and red = low performance, based on a linear scale of the best- to worst-result for each task.

the start of the paper. In this, we make reference to: (1) the best-performing method both in- and out-of-domain vs. the state-of-the-art (Tab. 3); (2) a heat map for each task indicating the convergence rate for each word representation, with and without updating (Fig. 1); (3) OOV accuracy both in-domain and out-of-domain for each task (Fig. 2); and (4) visualisation of the impact of updating on word embeddings, based on t-SNE (Fig. 3).

RQ1: Are word embeddings better than one-hot unigram features and Brown clusters? As shown in Tab. 3, the best-performing method for every task except in-domain Chunking is a word embedding method, although the precise method varies greatly. Fig. 1, on the other hand,

tells a more subtle story: the difference between UNIGRAM and the other word representations is relatively modest, esp. as the amount of training data increases. Additionally, the difference between BROWN and the word embedding methods is modest across all tasks. So, the overall answer would appear to be: yes for unigrams when there is little training data, but not really for BROWN.

RQ2: Do word embedding features require less training data? Fig. 1 shows that for POS-tagging and NER, with only several hundred training instances, word embedding features achieve superior results to UNIGRAM. For example, when trained with 561 instances, the POS-tagging model using SKIP-GRAM+UP embeddings is 5.3% above UNIGRAM; and when

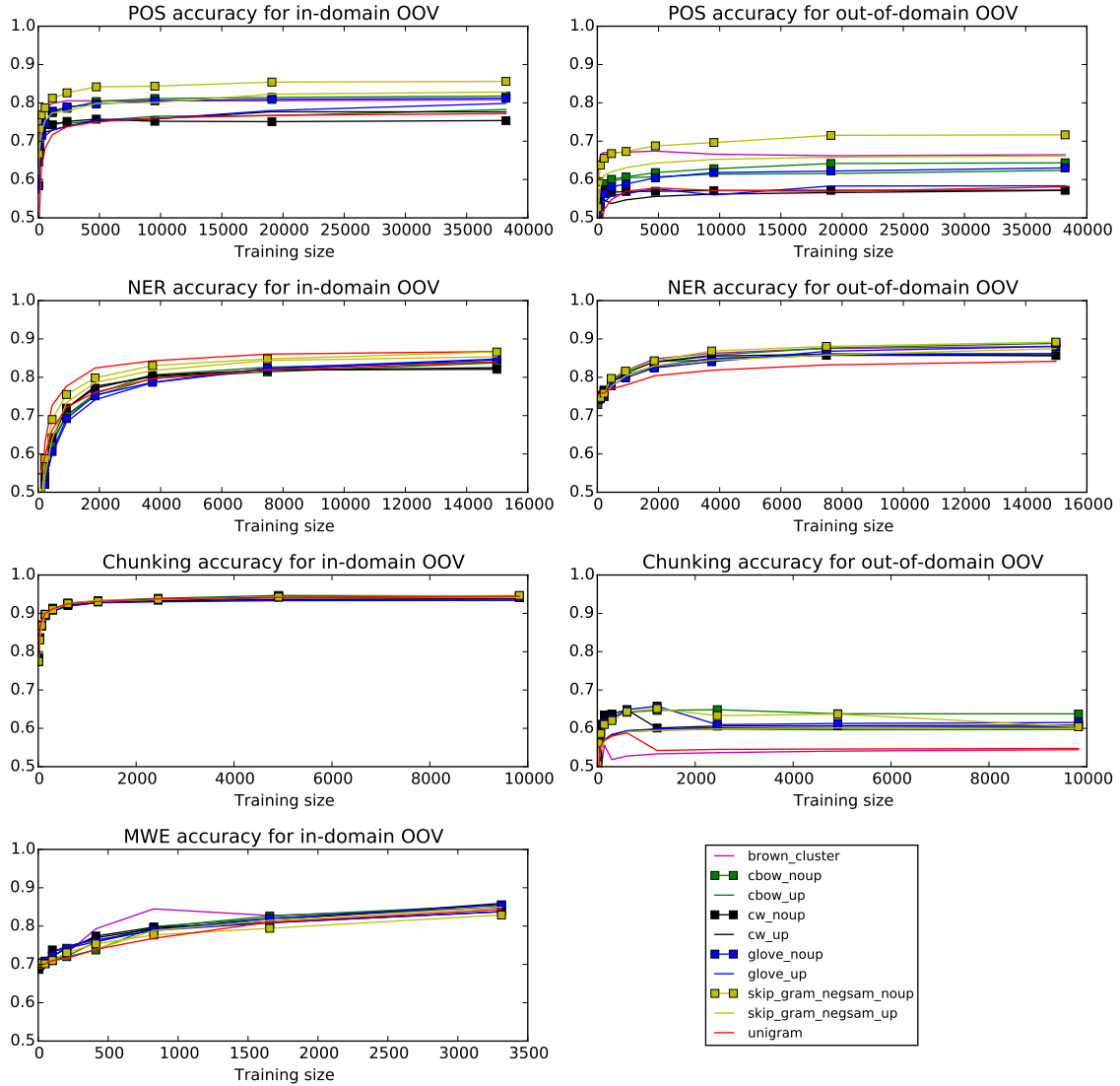


Figure 2: ACC over out-of-vocabulary (OOV) words for *in-domain* and *out-of-domain* test sets.

trained with 932 instances, the NER model using SKIP-GRAM is 11.7% above UNIGRAM. Similar improvements are also found for other types of word embeddings and BROWN, when the training set is small. However, all word representations perform similarly for Chunking regardless of training data size. For MWE, BROWN performs slightly better than the other methods when trained with approximately 25% of the training instances. Therefore, we conjecture that the POS-tagging and NER tasks benefit more from distributional similarity than Chunking and MWE.

RQ3: Does task-specific updating improve all word embeddings across all tasks? Based on Fig. 1, updating of word representations can equally correct poorly-learned word representations, and harm pre-trained representations, due to overfitting. For example, CW and GLOVE

perform significantly worse than SKIP-GRAM in both POS-tagging and NER without updating, but *with* updating, the gap between their results and the best performing method becomes smaller. In contrast, SKIP-GRAM performs worse over the test data with updating, despite the results on the development set improving by 1%.

To further investigate the effects of updating, we sampled 60 words and plotted the changes in their word embeddings under updating, using 2-d vector fields generated by using matplotlib and t-SNE (?). Half of the words were chosen manually to include known word clusters such as days of the week and names of countries; the other half were selected randomly. Additional plots with 100 randomly-sampled words and the top-100 most frequent words, for all the methods and all the tasks, can be found in the supplementary material and at <https://>

123abc123abd.wordpress.com/. In each plot, a single arrow signifies one word, pointing from the position of the original word embedding to the updated representation.

In Fig. 3, we show vector fields plots for Chunking and NER using SKIP-GRAM embeddings. For Chunking, most of the vectors were changed with similar magnitude, but in very different directions, including within the clusters of days of the week and country names. In contrast, for NER, there was more homogeneous change in word vectors belonging to the same cluster. This greater consistency is further evidence that semantic homogeneity appears to be more beneficial for NER than Chunking.

RQ4: What is the impact of word embeddings cross-domain and for OOV words? As shown in Tab. 3, results predictably drop when we evaluate out of domain. The difference is most pronounced for Chunking, where there is an absolute drop in F1 of around 30% for all methods, indicating that word embeddings and unigram features provide similar information for Chunking.

Another interesting observation is that updating often hurts out-of-domain performance because the distribution between domains is different. This suggests that, if the objective is to optimise performance across domains, it is best not to perform updating.

We also analyze performance on OOV words both in-domain and out-of-domain in Fig. 2. As expected, word embeddings and BROWN excel in out-of-domain OOV performance. Consistent with our overall observations about cross-domain generalisation, the OOV results are better when updating is not performed.

RQ5 Overall, are some word embeddings better than others? Comparing the different word embedding techniques over our four sequence labelling tasks, for the different evaluations (overall, out-of-domain and OOV), there is no clear winner among the word embeddings – for POS-tagging, SKIP-GRAM appears to have a slight advantage, but this does not generalise to other tasks.

While the aim of this paper was not to achieve the state of the art over the respective tasks, it is important to concede that our best (in-domain) results for NER, POS-tagging and Chunking are slightly worse than the state of the art (Tab. 3). The

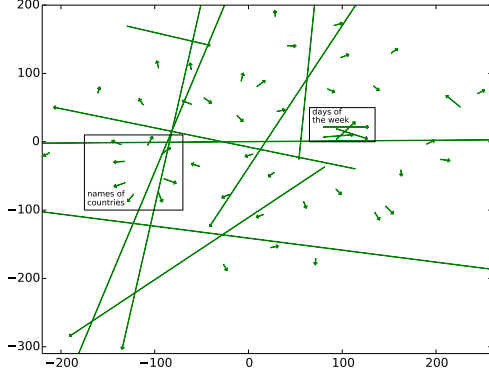
2.7% difference between our NER system and the best performing system is due to the fact that we use a first-order instead of a second-order CRF (?), and for the other tasks, there are similarly differences in the learner and the complexity of the features used. Another difference is that we tuned the hyperparameters with random search, to enable replication using the same random seed. In contrast, the hyperparameters for the state-of-the-art methods are tuned more extensively by experts, making them more difficult to reproduce.

5 Related Work

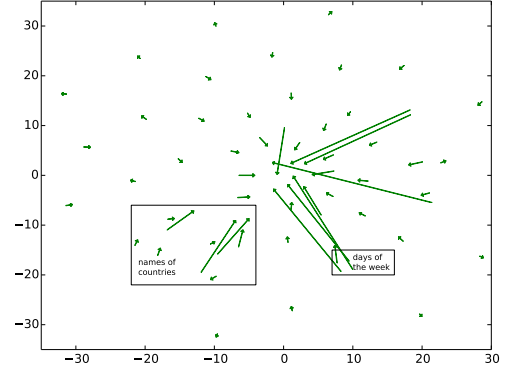
?) proposed a unified neural network framework that learns word embeddings and applied it for POS-tagging, Chunking, NER and semantic role labelling. When they combined word embeddings with hand crafted features (e.g., word suffixes for POS-tagging; gazetteers for NER) and applied other tricks like cascading and classifier combination, they achieved state-of-the-art performance. Similarly, ?) evaluated three different word representations on NER and Chunking, and concluded that unsupervised word representations improved NER and Chunking. They also found that combining different word representations can further improve performance. ?) also explored different ways of using word embeddings for NER. ?) and ?) found that BROWN clustering enhances Twitter POS tagging and MWE, respectively. Compared to previous work, we consider *more* word representations including the most recent work and evaluate them on *more* sequence labelling tasks, wherein the models are trained with training sets of varying size.

?) reported that direct use of word embeddings in dependency parsing did not show improvement. They achieved an improvement only when they performed hierarchical clustering of the word embeddings, and used features extracted from the cluster hierarchy. In a similar vein, ?) explored the use of word embeddings for constituency parsing and concluded that the information contained in word embeddings might duplicate the one acquired by a syntactic parser, unless the training set is extremely small. Other syntactic parsing studies that reported improvements by using word embeddings include ?), ?), ?) and ?).

Word embeddings have also been applied to other (non-sequential NLP) tasks like grammar induction (?), and semantic tasks such as seman-



(a) Chunking



(b) NER

Figure 3: A t-SNE plot of the impact of updating on SKIP-GRAM

tic relatedness, synonymy detection, concept categorisation, selectional preference learning and analogy (?).

6 Conclusions

We have performed an extensive extrinsic evaluation of five word embedding methods under fixed experimental conditions, and evaluated their applicability to four sequence labelling tasks: POS-tagging, Chunking, NER and MWE identification. We found that word embedding features reliably outperformed unigram features, especially with limited training data, but that there was relatively little difference over Brown clusters, and no one embedding method was consistently superior across the different tasks and settings. Word embeddings and Brown clusters were also found to improve out-of-domain performance and for OOV words. We expected a performance gap between the fixed and task-updated embeddings, but the observed difference was marginal. Indeed, we found that updating can result in overfitting. We also carried out preliminary analysis of the impact of updating on the vectors, a direction which we intend to pursue further.