



**Representing Meaning in Continuous Space: From  
Words to Sentences**

Felix Hill

This dissertation is submitted for the degree of Doctor of Philosophy



# Abstract

My abstract ...



# Declaration

This dissertation is my own work and contains nothing which is the outcome of work done in collaboration with others, except where specified in the text. This dissertation is not substantially the same as any that I have submitted for a degree or diploma or other qualification at any other university. This dissertation does not exceed the prescribed limit of 60 000 words.

Felix Hill  
April 2016



# Acknowledgements

My acknowledgements ...





# Contents

<b>1</b>	<b>Introduction</b>	<b>15</b>
<b>2</b>	<b>Neural Language Models</b>	<b>17</b>
<b>3</b>	<b>Distributed Representations of Words</b>	<b>19</b>
3.1	Learning Distributed Word Representations from Text . . . . .	19
3.2	Modelling Word Acquisition with Multi-Modal Data and Neural Language Models . . . . .	19
3.2.1	Introduction . . . . .	19
3.2.2	Model Design . . . . .	21
3.2.3	Information Sources . . . . .	24
3.2.4	Evaluation . . . . .	25
3.2.5	Results and Discussion . . . . .	26
3.2.6	Combining information sources . . . . .	27
3.2.7	Propagating input to abstract concepts . . . . .	28
3.2.8	Direct representation vs. propagation . . . . .	31
3.2.9	Source and quantity of perceptual input . . . . .	32
3.2.10	Conclusions . . . . .	32
3.3	Improving the Evaluation of Word Representations . . . . .	34
3.4	Sequence-to-Sequence Learning of Word Representations From Bilingual Data . . . . .	34
3.5	Introduction . . . . .	34
3.6	Learning Embeddings with Neural Language Models . . . . .	35
3.6.1	Monolingual Models . . . . .	35
3.6.2	Bilingual Representation-learning Models . . . . .	36
3.6.3	Neural Machine Translation Models . . . . .	37

3.7	Experiments . . . . .	38
3.7.1	Similarity and relatedness modelling . . . . .	39
3.7.2	Importance of training data quantity . . . . .	41
3.7.3	Analogy Resolution . . . . .	42
3.8	Effect of Target Language . . . . .	43
3.9	Overcoming the Vocabulary Size Problem . . . . .	45
3.10	How Similarity Emerges . . . . .	46
3.11	Conclusion . . . . .	48
<b>4</b>	<b>Learning to Represent Phrases</b>	<b>51</b>
4.1	Introduction . . . . .	51
4.2	Neural Language Model Architectures . . . . .	52
4.2.1	Long Short Term Memory . . . . .	53
4.2.2	Bag-of-Words NLMs . . . . .	54
4.2.3	Pre-trained Input Representations . . . . .	55
4.2.4	Training Objective . . . . .	55
4.2.5	Implementation Details . . . . .	55
4.3	Reverse Dictionaries . . . . .	56
4.3.1	Data Collection and Training . . . . .	57
4.3.2	Comparisons . . . . .	58
4.3.3	Reverse Dictionary Evaluation . . . . .	59
4.3.4	Results . . . . .	60
4.3.5	Qualitative Analysis . . . . .	61
4.3.6	Cross-Lingual Reverse Dictionaries . . . . .	63
4.3.7	Discussion . . . . .	64
4.4	General Knowledge (crossword) Question Answering . . . . .	64
4.4.1	Evaluation . . . . .	65
4.4.2	Benchmarks and Comparisons . . . . .	66
4.4.3	Results . . . . .	67
4.4.4	Qualitative Analysis . . . . .	68
4.5	Conclusion . . . . .	69
<b>5</b>	<b>Learning to Represent Sentences</b>	<b>71</b>
5.0.1	Introduction . . . . .	71
5.0.2	Distributed Sentence Representations . . . . .	72

5.0.3	Existing Models Trained on Text . . . . .	73
5.0.4	Models Trained on Structured Resources . . . . .	74
5.0.5	Novel Text-Based Models . . . . .	75
5.0.6	Training and Model Selection . . . . .	76
5.0.7	Evaluating Sentence Representations . . . . .	77
5.0.8	Supervised Evaluations . . . . .	78
5.0.9	Unsupervised Evaluations . . . . .	79
5.0.10	Results . . . . .	80
5.0.11	Discussion . . . . .	80
5.0.12	Conclusion . . . . .	83
<b>6</b>	<b>Representation Learning</b>	<b>85</b>
<b>7</b>	<b>Conclusion</b>	<b>87</b>
	<b>Bibliography</b>	<b>89</b>
<b>A</b>	<b>Extra Information</b>	<b>103</b>







# **Chapter 1**

## **Introduction**

Some introduction text





# **Chapter 2**

## **Neural Language Models**

Some NLM text



# Chapter 3

## Distributed Representations of Words

Some word embedding text

### 3.1 Learning Distributed Word Representations from Text

A brief history of how this was done...Bengio - Collobert - Turian - Mikolov

Contrast with 'predicting' models: LSA, Sahlgren, Baroni

Levy two results.

### 3.2 Modelling Word Acquisition with Multi-Modal Data and Neural Language Models

#### 3.2.1 Introduction

Multi-modal models that learn semantic representations from both language and information about the perceptible properties of concepts were originally motivated by parallels with human word learning [4] and evidence that many concepts are grounded in perception [9]. The perceptual information in such models is generally mined directly from images [38, 19] or from data collected in psychological studies [106, 100].

By exploiting the additional information encoded in perceptual input, multi-modal models can outperform language-only models on a range of semantic NLP tasks, including modelling similarity [20] and free association [106], predicting compositional-

ity [100] and concept categorization [104]. However, to date, this superiority has only been established when evaluating on concrete words such as *cat* or *dog*, rather than abstract concepts, such as *curiosity* or *loyalty*. Indeed, differences between abstract and concrete processing and representation [94, 53] suggest that conclusions about concrete concept learning may not necessarily hold in the general case. In this paper, we therefore focus on multi-modal models for learning both abstract and concrete concepts.

Although concrete concepts might seem more basic or fundamental, the vast majority of open-class, meaning-bearing words in everyday language are in fact abstract. 72% of the noun or verb tokens in the British National Corpus [71] are rated by human judges<sup>1</sup> as more abstract than the noun *war*, for instance, a concept many would already consider to be quite abstract. Moreover, abstract concepts by definition encode higher-level (more general) principles than concrete concepts, which typically reside naturally in a single semantic category or domain [30]. It is therefore likely that abstract representations may prove highly applicable for multi-task, multi-domain or transfer learning models, which aim to acquire ‘general-purpose’ conceptual knowledge without reference to a specific objective or task [28, 80].

Motivated by these observations, we introduce an architecture for learning both abstract and concrete representations that generalizes the skipgram model of [81] from corpus-based to multi-modal learning. The extended model is designed to reflect aspects of human word learning, in that it introduces more perceptual information about commonly-occurring concrete concepts and less information about rarer concepts.

We train our model on running-text language and two sources of perceptual descriptors for concrete nouns: the ESPGame dataset of annotated images [116] and the CSLB set of concept property norms [32]. We find that our model *combines* information from the different modalities more effectively than previous methods, resulting in an improved ability to model the USF free association gold standard [92] for concrete nouns. In addition, the architecture *propagates* the extra-linguistic input for concrete nouns to improve representations of abstract concepts more effectively than alternative methods. While this propagation can effectively extend the advantage of the multi-modal approach to many more concepts than simple concrete nouns, we observe that the benefit of adding perceptual input appears to decrease as target concepts become more abstract. Indeed, for the most abstract concepts of all, language-only models still

---

<sup>1</sup>Contributors to the USF dataset [92]

provide the most effective learning mechanism.

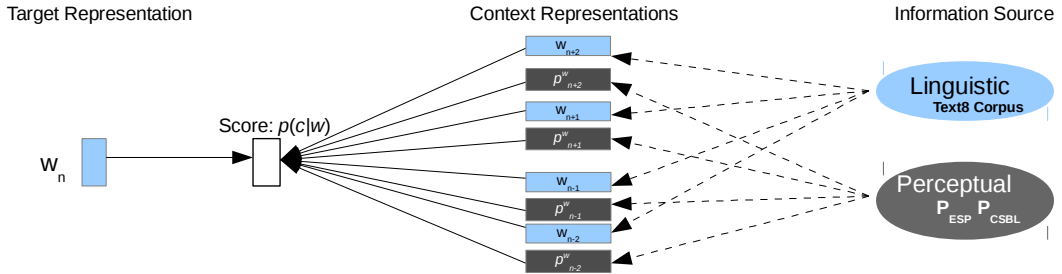
Finally, we investigate the optimum quantity and type of perceptual input for such models. Between the most concrete concepts, which can be effectively represented directly in the perceptual modality, and the most abstract concepts, which cannot, we identify a set of concepts that cannot be represented effectively directly in the perceptual modality, but still benefit from perceptual input propagated in the model via concrete concepts.

The motivation in designing our model and experiments is both practical and theoretical. Taken together, the empirical observations we present are potentially important for optimizing the learning of representations of concrete and abstract concepts in multi-modal models. In addition, they offer a degree of insight into the poorly understood issue of how abstract concepts may be encoded in human memory.

### 3.2.2 Model Design

Before describing how our multi-modal architecture encodes and integrates perceptual information, we first describe the underlying corpus-based representation learning model.

**Language-only Model** Our multi-modal architecture builds on the continuous log-linear skip-gram language model proposed by [81]. This model learns lexical representations in a similar way to neural-probabilistic language models (NPLM) but without a non-linear hidden layer, a simplification that facilitates the efficient learning of large vocabularies of dense representations, generally referred to as *embeddings* [112]. Embeddings learned by the model achieve state-of-the-art performance on several evaluations including sentence completion and analogy modelling [81].



**Figure 3.1:** Our multi-modal model architecture. Light boxes are elements of the original [81] model. For target words  $w_n$  in the domain of  $\mathbf{P}$ , the model updates based on corpus context words  $w_{n+i}$  then on words  $p_{n+i}^w$  in perceptual psuedo-sentences. Otherwise, updates are based solely on the  $w_{n+i}$ .

For each word type  $w$  in the vocabulary  $V$ , the model learns both a ‘target-embedding’  $r_w \in \mathbb{R}^d$  and a ‘context-embedding’  $\hat{r}_w \in \mathbb{R}^d$  such that, given a target word, its ability to predict nearby context words is maximized. The probability of seeing context word  $c$  given target  $w$  is defined as:

$$p(c|w) = \frac{e^{\hat{r}_c \cdot r_w}}{\sum_{v \in V} e^{\hat{r}_v \cdot r_w}}$$

The model learns from a set of target-word, context-word pairs, extracted from a corpus of sentences as follows. In a given sentence  $S$  (of length  $N$ ), for each position  $n \leq N$ , each word  $w_n$  is treated in turn as a target word. An integer  $t(n)$  is then sampled from a uniform distribution on  $\{1, \dots, k\}$ , where  $k > 0$  is a predefined maximum context-window parameter. The pair tokens  $\{(w_n, w_{n+j}) : -t(n) \leq j \leq t(n), w_i \in S\}$  are then appended to the training data. Thus, target/context training pairs are such that (i) only words within a  $k$ -window of the target are selected as context words for that target, and (ii) words closer to the target are more likely to be selected than those further away.

The training objective is then to maximize the log probability  $T$  across of all such examples from  $S$ , and then across all sentences in the corpus:

$$T = \frac{1}{N} \sum_{n=1}^N \sum_{-t(n) \leq j \leq t(n), j \neq 0} \log(p(w_{n+j}|w_n))$$

The model free parameters (target-embeddings and context-embeddings of dimension  $d$  for each word in the corpus with frequency above a certain threshold  $f$ ) are updated according to stochastic gradient descent and backpropation, with learning rate controlled by Adagrad [35]. For efficiency, the output layer is encoded as a hierarchical softmax function based on a binary Huffman tree [90].

As with other distributional architectures, the model captures conceptual semantics by exploiting the fact that words appearing in similar linguistic contexts are likely to have similar meanings. Informally, the model adjusts its embeddings to increase the ‘probability’ of seeing the language in the training corpus. Since this probability increases with the  $p(c|w)$ , and the  $p(c|w)$  increase with the dot product  $\hat{r}_v \cdot r_c$ , the updates have the effect of moving each target-embedding incrementally ‘closer’ to the context-embeddings of its collocates. In the target-embedding space, this results in embeddings of concept words that regularly occur in similar contexts moving closer together.

**Multi-modal Extension** We extend the [81] architecture via a simple means of introducing perceptual information that aligns with human language learning. Based on the assumption that frequency in domain-general linguistic corpora correlates with the likelihood of ‘experiencing’ a concept in the world [21, 23], perceptual information is introduced to the model whenever designated concrete concepts are encountered in the running-text linguistic input. This has the effect of introducing more perceptual input for commonly experienced concrete concepts and less input for rarer concrete concepts.

To implement this process, perceptual information is extracted from external sources and encoded in an associative array  $\mathbf{P}$ , which maps (typically concrete) words  $w$  to bags of perceptual features  $\mathbf{b}(w)$ . The construction of this array depends on the perceptual information source; the process for our chosen sources is detailed in Section 3.2.3.

Training our model begins as before on running-text. When a sentence  $S_m$  containing a word  $w$  in the domain of  $\mathbf{P}$  is encountered, the model completes training on  $S_m$  and begins learning from a perceptual pseudo-sentence  $\hat{S}(w)$ .  $\hat{S}_m(w)$  is constructed by randomly sampling features from  $\mathbf{b}(w)$  to occupy positions before and instances of  $w$ , so that  $\hat{S}_m(w)$  is the same length as  $S_m$  (see Figure 3.2). Once training on  $\hat{S}_m(w)$  is completed, the model reverts to the next ‘real’ (linguistic) sentence  $S_{m+1}$ , and the process continues. Thus, when a concrete concept is encountered in the corpus, its embedding is first updated based on language (moved incrementally closer to concepts appearing in similar linguistic contexts), and then on perception (moved incrementally closer to concepts with the same or similar perceptual features).

For greater flexibility, we introduce a parameter  $\alpha$  reflecting the raw quantity of perceptual information relative to linguistic input. When  $\alpha = 2$ , two pseudo-sentences are generated and inserted for every corpus occurrence of a token from the domain of  $\mathbf{P}$ . For non-integral  $\alpha$ , the number of sentences inserted is  $\lfloor \alpha \rfloor$ , and a further sentence is added with probability  $\alpha - \lfloor \alpha \rfloor$ .

In all experiments reported in the following sections we set the window size parameter  $k = 5$  and the minimum frequency parameter  $f = 3$ , which guarantees that the model learns embeddings for all concepts in our evaluation sets. While the model learns both target and context-embeddings for each word in the vocabulary, we conduct our experiments with the target embeddings only. We set the dimension parameter  $d = 300$  as this produces high quality embeddings in the language-only case [81].

$\hat{S}(\textit{crocodile}) = \textbf{Crocodile}$  legs **crocodile** teeth **crocodile** teeth **crocodile** scales **crocodile** green **crocodile**.

$\hat{S}(\textit{screwdriver}) = \textbf{Screwdriver}$  handle **screwdriver** flat **screwdriver** long **screwdriver** handle **screwdriver** head.

**Figure 3.2:** Example pseudo-sentences generated by our model.

### 3.2.3 Information Sources

We construct the associative array of perceptual information  $\mathbf{P}$  from two sources typical of those typically used for multi-modal semantic models.

**ESPGame Dataset** The ESP-Game dataset (ESP) [116] consists of 100,000 images, each annotated with a list of lexical concepts that appear in that image.

For any concept  $w$  identified in an ESP image, we construct a corresponding bag of features  $\mathbf{b}(w)$ . For each ESP image  $I$  that contains  $w$ , we append the other concept tokens identified in  $I$  to  $\mathbf{b}(w)$ . Thus, the more frequently a concept co-occurs with  $w$  in images, the more its corresponding lexical token occurs in  $\mathbf{b}(w)$ . The array  $\mathbf{P}_{\text{ESP}}$  in this case then consists of the  $(w, \mathbf{b}(w))$  pairs.

**CSLB Property Norms** The Centre for Speech, Language and the Brain norms (CSLB) [32] is a recently-released dataset containing semantic properties for 638 concrete concepts produced by human annotators. The CSLB dataset was compiled in the same way as the [79] property norms used widely in multi-modal models [106, 100]; we use CSLB because it contains more concepts. For each concept, the proportion of the 30 annotators that produced a given feature can also be employed as a measure of the strength of that feature.

When encoding the CSLB data in  $\mathbf{P}$ , we first map properties to lexical forms (e.g. *is\_green* becomes *green*). By directly identifying perceptual features and linguistic forms in this way, we treat features observed in the perceptual data as (sub)concepts to be acquired via the same multi-modal input streams and stored in the same domain-general memory as the evaluation concepts. This design decision in fact corresponds to a view of cognition that is sometimes disputed [41]. In future studies we hope to compare the present approach to architectures with domain-specific conceptual memories.

For each concept  $w$  in CSLB, we then construct a feature bag  $\mathbf{b}(w)$  by append-



ESPGame		CSLB	
Image 1	Image 2	Crocodile	Screwdriver
red	wreck	has 4 legs (7)	has handle (28)
chihuaua	cyan	has tail (18)	has head (5)
eyes	man	has jaw (7)	is long (9)
little	crash	has scales (8)	is plastic (18)
ear	accident	has teeth (20)	is metal (28)
nose	street	is green (10)	
small		is large (10)	

**Table 3.1:** Concepts identified in images in the ESP Game (left) and features produced for concepts by human annotators in the CSLB dataset (with feature strength, max=30).

ing lexical forms to  $\mathbf{b}(w)$  such that the count of each feature word is equal to the strength of that feature for  $w$ . Thus, when features are sampled from  $\mathbf{b}(w)$  to create pseudo-sentences (as detailed previously) the probability of a feature word occurring in a sentence reflects feature strength. The array  $\mathbf{P}_{\text{CSLB}}$  then consists of all  $(w, \mathbf{b}(w))$  pairs.

**Linguistic Input** The linguistic input to all models is the 400m word Text8 Corpus<sup>2</sup> of Wikipedia text, split into sentences and with punctuation removed.

### 3.2.4 Evaluation

We evaluate the quality of representations by how well they reflect *free association* scores, an empirical measure of cognitive conceptual proximity. The University of South Florida Norms (USF) [92] contain free association scores for over 40,000 concept pairs, and have been widely used in NLP to evaluate semantic representations [4, 38, 106, 100]. Each concept that we extract from the USF database has also been rated for conceptual concreteness on a Likert scale of 1-7 by at least 10 human annotators. Following previous studies [57, 106], we measure the (Spearman  $\rho$ ) correlation between association scores and the cosine similarity of vector representations.

We create separate abstract and concrete concept lists by ranking the USF concepts according to concreteness and sampling at random from the first and fourth quartiles. We also introduce a complementary noun/verb dichotomy,<sup>3</sup> on the intu-

<sup>2</sup>From <http://mattmahoney.net/dc/textdata.html>

<sup>3</sup>Based on the majority POS-tag of words in the lemmatized British National Corpus [71]

Concept 1	Concept 2	Assoc.
abdomen (6.83)	stomach (6.04)	0.566
throw (4.05)	ball (6.08)	0.234
hope (1.18)	glory (3.53)	0.192
egg (5.79)	milk (6.66)	0.012

**Table 3.2:** Example concept pairs (with mean concreteness rating) and free-association scores from the USF dataset.

Concept Type	List	Pairs	Examples
concrete nouns	541	1418	<i>yacht, cup</i>
abstract nouns	100	295	<i>fear, respect</i>
all nouns	666	1815	<i>fear, cup</i>
concrete verbs	50	66	<i>kiss, launch</i>
abstract verbs	50	127	<i>differ, obey</i>
all verbs	100	221	<i>kiss, obey</i>

**Table 3.3:** Details the subsets of USF data used in our evaluations, downloadable from our website.

ition that information propagation may occur differently from noun to noun or from noun to verb (because of their distinct structural relationships in sentences). The abstract/concrete and noun/verb dichotomies yield four distinct concept lists. For consistency, the concrete noun list is filtered so that all concrete noun concepts  $w$  have perceptual representations  $\mathbf{b}(w)$  in both  $\mathbf{P}_{\text{ESP}}$  and  $\mathbf{P}_{\text{CSLB}}$ . For each of the four resulting concept lists  $C$  (concrete/abstract, noun/verb), a corresponding set of evaluation pairs  $\{(w_1, w_2) \in \text{USF} : w_1, w_2 \in C\}$  is extracted (see Table 3 for details).

### 3.2.5 Results and Discussion

Our experiments were designed to answer four questions, outlined in the following subsections: (1) Which model architectures perform best at *combining* information pertinent to multiple modalities when such information exists explicitly (as common for concrete concepts)? (2) Which model architectures best propagate perceptual information to concepts for which it does not exist explicitly (as is common for abstract concepts)? (3) Is it preferable to include all of the perceptual input that can be obtained from a given source, or to filter this input stream in some way? (4) How much perceptual vs. linguistic input is optimal for learning various concept types?

### 3.2.6 Combining information sources

To evaluate our approach as a method of information combination we compared its performance on the concrete noun evaluation set against alternative methods. When implementing the alternatives, we first encoded the perceptual input directly into sparse feature vectors, with coordinates for each of the 2726 features in CSLB and for each of the 100,000 images in ESP.

The first alternative is simple concatenation of these perceptual vectors with linguistic vectors embeddings learned by the [81] model on the Text8 Corpus. In the second alternative, proposed for multi-modal models by [106], *canonical correlation analysis* (CCA) [47] was applied to the vectors of both modalities. This yields reduced-dimensionality representations that preserve underlying inter-modal correlations, which are then concatenated. The final alternative, proposed by [20] involves applying Singular Value Decomposition (SVD) to the matrix of concatenated multi-modal representations, yielding smoothed representations.<sup>4</sup>

We compare these alternatives to our proposed model with  $\alpha = 1$ . In The CSLB and ESP models, all training pseudo-sentences are generated from the arrays  $\mathbf{P}_{\text{CSLB}}$  and  $\mathbf{P}_{\text{ESP}}$  respectively. In the models classed as *CSLB&ESP*, a random choice between  $\mathbf{P}_{\text{CSLB}}$  and  $\mathbf{P}_{\text{ESP}}$  is made every time perceptual input is included (so that the overall quantity of perceptual information is the same).

As shown in Figure 3.3 (left side), the embeddings learned by our model achieve a higher correlation with the USF data than simple concatenation, CCA and SVD regardless of perceptual input source. With the optimal perceptual source (ESP only), for instance, the correlation is 11% higher than the next best alternative method, CCA.

One possible factor behind this improvement is that, in our model, the learned representations fully integrate the two modalities, whereas for both CCA and the concatenation method each representation feature (whether of reduced dimension or not) corresponds to a particular modality. This deeper integration may help our architecture to overcome the challenges inherent in information combination such as inter-modality differences in information content and representation sparsity.

---

<sup>4</sup>CCA was implemented using the *CCA* package in R. SVD was implemented using the Python *sparsesvd* package, with truncation factor  $k = 1024$  as per [20].

### 3.2.7 Propagating input to abstract concepts

To test the process of information propagation in our model, we evaluated the learned embeddings of more abstract concepts. We compared our approach with two recently-proposed alternative methods for inferring perceptual features when explicit perceptual information is unavailable.

**Johns and Jones** In the method of [62], pseudo-perceptual representations for target concepts without a perceptual representations (uni-modal concepts) are inferred as a weighted average of the perceptual representations of concepts that do have such a representation (bi-modal concepts).

In the first step of their two-step method, for each uni-modal concept  $\mathbf{k}$ , a quasi-perceptual representation is computed as an average of the perceptual representations of bi-modal concepts, weighted by the proximity between each of these concepts and  $\mathbf{k}$

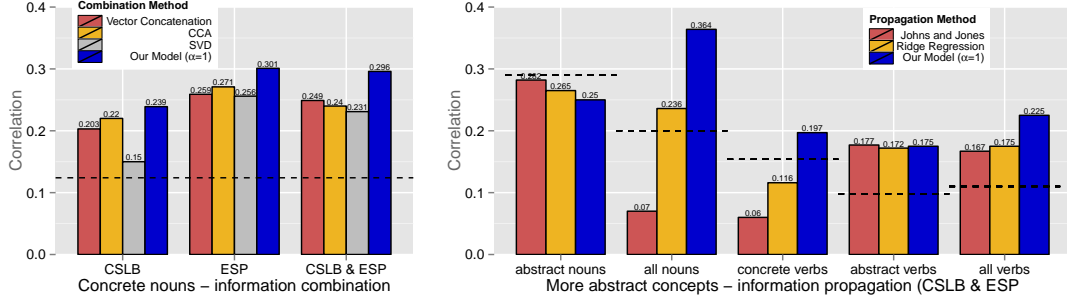
$$\mathbf{k}^p = \sum_{\mathbf{c} \in \bar{C}} S(\mathbf{k}^l, \mathbf{c}^l)^\lambda \cdot \mathbf{c}^p$$

where  $\bar{C}$  is the set of bi-modal concepts,  $\mathbf{c}^p$  and  $\mathbf{k}^p$  are the perceptual representations for  $\mathbf{c}$  and  $\mathbf{k}$  respectively, and  $\mathbf{c}^l$  and  $\mathbf{k}^l$  the linguistic representations. The exponent parameter  $\lambda$  reflects the learning rate.

In step two, the initial quasi-perceptual representations are inferred for a second time, but with the weighted average calculated over the perceptual or initial quasi-perceptual representations of all other words, not just those that were originally bi-modal. As with [62], we set the learning rate parameter  $\lambda$  to be 3 in the first step and 13 in the second.

**Ridge Regression** An alternative, proposed for the present purpose by [ref. withdrawn for review], uses *ridge regression* [91]. Ridge regression is a variant of least squares regression in which a regularization term is added to the training objective to favor solutions with certain properties.

For bi-modal concepts of dimension  $n_p$ , we apply ridge regression to learn  $n_p$  linear functions  $f_i : \mathbb{R}^{n_l} \rightarrow \mathbb{R}$  that map the linguistic representations (of dimension  $n_l$ ) to a particular perceptual feature  $i$ . These functions are then applied together to map the linguistic representations of uni-modal concepts to full quasi-perceptual representations.



**Figure 3.3:** The proposed approach compared with other methods of information combination (left) and propagation. Dashed lines indicate language-only model baseline.

Following [ref. withdrawn for review], we take the Euclidian  $l_2$  norm of the inferred parameter vector as the regularization term. This ensures that the regression favors lower coefficients and a smoother solution function, which should provide better generalization performance than simple linear regression. The objective for learning the  $f_i$  is then to minimize

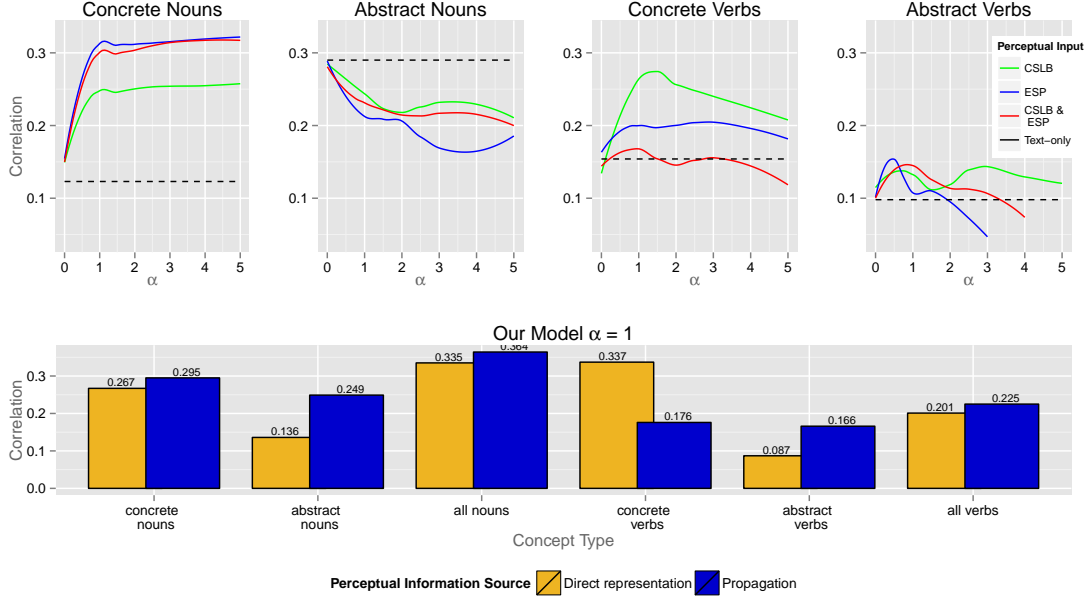
$$\|\mathbf{a}X - Y_i\|_2^2 + \|\mathbf{a}\|_2^2$$

where  $\mathbf{a}$  is the vector of regression coefficients,  $X$  is a matrix of linguistic representations and  $Y_i$  a vector of the perceptual feature  $i$  for the set of bi-modal concepts.

**Comparisons** We applied the Johns and Jones method and ridge regression starting from linguistic embeddings acquired by the [81] model on the Text8 Corpus, and concatenated the resulting pseudo-perceptual and linguistic representations. As with the implementation of our model, the perceptual input for these alternative models was limited to concrete nouns (i.e. concrete nouns were the only bi-modal concepts in the models).

Figure 3.3 (right side) illustrates the propagation performance of the three models. While the correlations overall may seem somewhat low, this is a consequence of the difficulty of modeling the USF data. In fact, the performance of both the language-only model and our multi-modal extension across the concept types, ranging from .18–.36, is equal to or higher than equivalent models evaluated on the same data previously [38, 106, 105].

For learning representations of concrete verbs, our approach achieves a 69% increase in performance over the next best alternative. The performance of the model on abstract verbs is marginally inferior to Johns and Jones’ method. Nevertheless, the clear advantage for concrete verbs makes our model the best choice for learning rep-



**Figure 3.4:** Top: Comparing the strategy of directly representing abstract concepts from perceptual information where available (yellow bars) vs. propagating via concrete concepts. Bottom: The effect of increasing  $\alpha$  on correlation with USF pairs (Spearman  $\rho$ ) for each concept type. Horizontal dashed lines indicate language-only model baseline.

representations of verbs in general, as shown by performance on the set *all verbs*, which also includes mixed abstract-concrete pairs.

Our model is also marginally inferior to alternative approaches in learning representations of abstract nouns. However, in this case, no method improves on the linguistic-only baseline. It is possible that perceptual information is simply so removed from the core semantics of these concepts that they are best acquired via the linguistic medium alone, regardless of learning mechanism. The moderately inferior performance of our method in such cases is likely caused by its greater inherent inter-modal dependence compared with methods that simply concatenate uni-modal representations. When the perceptual signal is of low quality, this greater inter-modal dependence allows the linguistic signal to be obscured. The trade-off, however, is the higher quality joint representations when the perceptual signal is of higher-quality, exemplified by the fact that our proposed approach outperforms alternatives on the set *all nouns*, which includes the more concrete nouns.

### 3.2.8 Direct representation vs. propagation

Although property norm datasets such as the CSLB data typically consist of perceptual feature information for concrete nouns only, image-based datasets such as ESP do contain information on more abstract concepts, which was omitted from the previous experiments. Indeed, image banks such as Google Images contain millions of photographs portraying quite abstract concepts, such as *love* or *war*. On the other hand, encodings or descriptions of abstract concepts are generally more subjective and less reliable than those of concrete concepts [122]. We therefore investigated whether or not it is preferable to include this additional information as model input or to restrict perceptual input to concrete nouns as previously.

Of our evaluation sets, it was possible to construct from ESP (and add to  $\mathbf{P}_{\text{ESP}}$ ) representations for all of the concrete verbs, and for approximately half of the abstract verbs and abstract nouns. Figure 3.4 (top), shows the performance of a our model trained on all available perceptual input versus the model in which the perceptual input was restricted to concrete nouns.

The results reflect a clear manifestation of the abstract/concrete distinction. Concrete verbs behave similarly to concrete nouns, in that they can be effectively represented directly from perceptual information sources. The information encoded in these representations is beneficial to the model and increases performance. In contrast, constructing ‘perceptual’ representations of abstract verbs and abstract nouns directly from perceptual information sources is clearly counter-productive (to the extent that performance also degrades on the combined sets *all nouns* and *all verbs*). It appears in these cases that the perceptual input acts to obscure or contradict the otherwise useful signal inferred from the corpus.

As shown in the previous section, the inclusion of any form of perceptual input inhibits the learning of abstract nouns. However, this is not the case for abstract verbs. Our model learns higher quality representations of abstract verbs when perceptual input is restricted to concrete nouns than when no perceptual input is included whatsoever *and* when perceptual input is included for both concrete nouns and abstract verbs. This supports the idea of a gradual scale of concreteness: the most concrete concepts can be effectively represented directly in the perceptual modality; somewhat more abstract concepts cannot be represented directly in the perceptual modality, but have representations that are improved by propagating perceptual input from concrete concepts via language; and the most abstract concepts are best acquired via language

alone.

### 3.2.9 Source and quantity of perceptual input

For different concept types, we tested the effect of varying the proportion of perceptual to linguistic input (the parameter  $\alpha$ ). Perceptual input was restricted to concrete nouns as in Sections 3.1-3.2.

As shown in Figure 3.4, performance on concrete nouns improves (albeit to a decreasing degree) as  $\alpha$  increases. When learning concrete noun representations, linguistic input is apparently redundant if perceptual input is of sufficient quality and quantity. For the other concept types, in each case there is an optimal value for  $\alpha$  in the range .5–2, above which perceptual input obscures the linguistic signal and performance degrades. The proximity of these optima to 1 suggests that for optimal learning, when a concrete concept is experienced approximately equal weight should be given to available perceptual and linguistic information.

### 3.2.10 Conclusions

Motivated by the notable prevalence of abstract concepts in everyday language, and their likely importance to flexible, general-purpose representation learning, we have investigated how abstract and concrete representations can be acquired by multi-modal models. In doing so, we presented a simple and easy-to-implement architecture for acquiring semantic representations of both types of concept from linguistic and perceptual input.

While neuro-probabilistic models have been applied to the problem of multi-modal representation learning previously [108, 123] our model and experiments develop this work in several important ways. First, we address the problem of learning abstract concepts. By isolating concepts of different concreteness and part-of-speech in our evaluation sets, and separating the processes of information combination and propagation, we demonstrate that the multi-modal approach is indeed effective for some, but perhaps not all, abstract concepts. In addition, our model introduces a clear parallel with human language learning. Perceptual input is introduced precisely when concrete concepts are ‘experienced’ by the model in the corpus text, much like a language learner experiencing concrete entities via sensory perception.

Taken together, our findings indicate the utility of distinguishing three concept



types when learning representations in the multi-modal setting.

**Type I** Concepts that can be effectively represented directly in the perceptual modality. For such concepts, generally concrete nouns or concrete verbs, our proposed approach provides a simple means of combining perceptual and linguistic input. The resulting multi-modal representations are of higher quality than those learned via other approaches, resulting in a performance improvement of over 10% in modelling free association.

**Type II** Concepts, including abstract verbs, that cannot be effectively represented directly in the perceptual modality, but whose representations can be improved by joint learning from linguistic input and perceptual information about related concepts. Our model can effectively propagate perceptual input (exploiting the relations inferred from the linguistic input) from Type I concepts to enhance the representations of Type II concepts above the language-only baseline. Because of the frequency of abstract concepts, such propagation extends the benefit of the multi-modal approach to a far wider range of language than models based solely in the concrete domain.

**Type III** Concepts, such as abstract nouns, which are more effectively learned via language-only models than multi-modal models. Neither the model we introduce here nor other proposed propagation methods achieve an improvement in representation quality for these concepts over the language-only baseline. Of course, it is an empirical question whether a multi-modal approach could ever enhance the representation learning of these concepts, one with potential implications for cognitive theories of grounding (a topic of much debate in psychology [44, 8]).

Additionally, we investigated the optimum type and quantity of perceptual input for learning concepts of different types. We showed that too much perceptual input can result in degraded representations. For concepts of type I and II, the optimal quantity resulted from setting  $\alpha = 1$ ; i.e. whenever a concrete concept was encountered, the model learned from an equal number of language-based and perception-based examples. While we make no formal claims here, such observations may ultimately provide insight into human language learning and semantic memory.

In future we will address the question of whether Type III concepts can ever be enhanced via multi-modal learning, and investigate multi-modal models that optimally learn concepts of each type. This may involve filtering the perceptual input stream for

concepts according to concreteness, and possibly more elaborate model architectures that facilitate distinct representational frameworks for abstract and concrete concepts.

### **3.3 Improving the Evaluation of Word Representations**

### **3.4 Sequence-to-Sequence Learning of Word Representations From Bilingual Data**

### **3.5 Introduction**

It is well known that word representations can be learned from the distributional patterns in corpora. Originally, such representations were constructed by counting word co-occurrences, so that the features in one word’s representation corresponded to other words [69, 113]. Neural language models, an alternative method for learning word representations, use language data to optimise (latent) features with respect to a language modelling objective. The objective can be to predict either the next word given the initial words of a sentence [10, 88, 28], or simply a nearby word given a single cue word [84, 97]. The representations learned by neural models (sometimes called *embeddings*) perform very effectively when applied as pre-trained features in a range of NLP applications and tasks [7].

Despite these clear results, it is not well understood how the architecture of neural models affects the information encoded in their embeddings. Here we contribute to this understanding by considering the embeddings learned by architectures with a very different objective function: *neural machine translation* (NMT) *models*. NMT models have recently emerged as an alternative to statistical, phrase-based translation models, and are beginning to achieve impressive translation performance [63, 33, 110].

We show that NMT models are not only a potential new direction for machine translation, but are also an effective means of learning word embeddings. Specifically, translation-based embeddings encode information relating to conceptual similarity (rather than non-specific relatedness or association) and lexical syntactic role more effectively than embeddings from monolingual neural language models. We demonstrate that these properties persist when translating between different language pairs (English-French and English-German). Further, based on the observation of subtle

language-specific effects in the embedding spaces, we conjecture as to why similarity dominates over other semantic relations in translation embedding spaces. Finally, we discuss a potential limitation of the application of NMT models for embedding learning - the computational cost of training large vocabularies of embeddings - and show that a novel method for overcoming this issue preserves the aforementioned properties of translation-based embeddings.

## 3.6 Learning Embeddings with Neural Language Models

All neural language models, including NMT models, learn real-valued embeddings (of specified dimension) for words in some pre-specified vocabulary,  $V$ , covering many or all words in their training corpus. At each training step, a ‘score’ for the current training example (or batch) is computed based on the embeddings in their current state. This score is compared to the model’s objective function, and the error is backpropagated to update both the model weights (affecting how the score is computed from the embeddings) and the embedding features themselves. At the end of this process, the embeddings should encode information that enables the model to optimally satisfy its objective.

### 3.6.1 Monolingual Models

In the original neural language model [10] and subsequent variants [28], training examples consist of an ordered sequence of  $n$  words, with the model trained to predict the  $n$ -th word given the first  $n - 1$  words. The model first represents the input as an ordered sequence of embeddings, which it transforms into a single fixed length ‘hidden’ representation, generally by concatenation and non-linear projection. Based on this representation, a probability distribution is computed over the vocabulary, from which the model can sample a guess for the next word. The model weights and embeddings are updated to maximise the probability of correct guesses for all sentences in the training corpus.

More recent work has shown that high quality word embeddings can be learned via simpler models with no nonlinear hidden layer [84, 97]. Given a single word or unordered window of words in the corpus, these models predict which words will oc-

cur nearby. For each word  $w$  in  $V$ , a list of training cases  $(w, c) : c \in V$  is extracted from the training corpus according to some algorithm. For instance, in the *skipgram* approach [84], for each ‘cue word’  $w$  the ‘context words’  $c$  are sampled from windows either side of tokens of  $w$  in the corpus (with  $c$  more likely to be sampled if it occurs closer to  $w$ ).<sup>5</sup> For each  $w$  in  $V$ , the model initialises both a cue-embedding, representing the  $w$  when it occurs as a cue-word, and a context-embedding, used when  $w$  occurs as a context-word. For a cue word  $w$ , the model uses the corresponding cue-embedding and all context-embeddings to compute a probability distribution over  $V$  that reflects the probability of a word occurring in the context of  $w$ . When a training example  $(w, c)$  is observed, the model updates both the cue-word embedding of  $w$  and the context-word embeddings in order to increase the conditional probability of  $c$ .

### 3.6.2 Bilingual Representation-learning Models

Various studies have demonstrated that word representations can also be effectively learned from bilingual corpora, aligned at the document, paragraph or word level [46, 118, 83, 50, 22]. These approaches aim to represent the words from two (or more) languages in a common vector space so that words in one language are close to words with similar or related meanings in the other. The resulting multilingual embedding spaces have been effectively applied to bilingual lexicon extraction [46, 118, 83] and document classification [67, 50, 22, 68].

We focus our analysis on two representatives of this class of (non-NMT) bilingual model. The first is that of [50], whose embeddings improve on the performance of [67] in document classification applications. As with the NMT models introduced in the next section, this model can be trained directly on bitexts aligned only at the sentence rather than word level. When training, for aligned sentences  $S_E$  and  $S_F$  in different languages, the model computes representations  $R_E$  and  $R_F$  by summing the embeddings of the words in  $S_E$  and  $S_F$  respectively. The embeddings are then updated to minimise the divergence between  $R_E$  and  $R_F$  (since they convey a common meaning). A noise-contrastive loss function ensures that the model does not arrive at trivial (e.g. all zero) solutions to this objective. [50] show that, despite the lack of prespecified word alignments, words in the two languages with similar meanings converge in the bilingual embedding space.<sup>6</sup>

---

<sup>5</sup> Subsequent variants use different algorithms for selecting the  $(w, c)$  from the training corpus [52, 72]

<sup>6</sup>The models of [22] and [50] both aim to minimise the divergence between source and target lan-

The second model we examine is that of [37]. Unlike the models described above, [37] showed explicitly that projecting word embeddings from two languages (learned independently) into a common vector space can favourably influence the orientation of word embeddings when considered in their monolingual subspace; i.e relative to other words in their own language. In contrast to the other models considered in this paper, the approach of [37] requires bilingual data to be aligned at the word level.

### 3.6.3 Neural Machine Translation Models

The objective of NMT is to generate an appropriate sentence in a target language  $S_t$  given a sentence  $S_s$  in the source language (see e.g. [63, 110]). As a by-product of learning to meet this objective, NMT models learn distinct sets of embeddings for the vocabularies  $V_s$  and  $V_t$  in the source and target languages respectively.

Observing a training case  $(S_s, S_t)$ , these models represent  $S_s$  as an ordered sequence of embeddings of words from  $V_s$ . The sequence for  $S_s$  is then encoded into a single representation  $R_s$ .<sup>7</sup> Finally, by referencing the embeddings in  $V_t$ ,  $R_s$  and a representation of what has been generated thus far, the model decodes a sentence in the target language word by word. If at any stage the decoded word does not match the corresponding word in the training target  $S_t$ , the error is recorded. The weights and embeddings in the model, which together parameterise the encoding and decoding process, are updated based on the accumulated error once the sentence decoding is complete.

Although NMT models can differ in their low-level architecture [63, 26, 5], the translation objective exerts similar pressure on the embeddings in all cases. The source language embeddings must be such that the model can combine them to form single representations for ordered sequences of multiple words (which in turn must enable the decoding process). The target language embeddings must facilitate the process of decoding these representations into correct target-language sentences.

---

guage sentences represented as sums of word embeddings. Because of these similarities, we do not compare with both in this paper.

<sup>7</sup>Alternatively, subsequences (phrases) of  $S_s$  may be encoded at this stage in place of the whole sentence [5].

### 3.7 Experiments

To learn translation-based embeddings, we trained two different NMT models. The first is the RNN encoder-decoder, *RNNenc* [26], which uses a recurrent-neural-network to encode all of the source sentence into a single vector on which the decoding process is conditioned. The second is the *RNN Search* architecture [5], which was designed to overcome limitations exhibited by the RNN encoder-decoder when translating very long sentences. *RNN Search* includes a *attention* mechanism, an additional feed-forward network that learns to attend to different parts of the source sentence when decoding each word in the target sentence.<sup>8</sup> Both models were trained on a 348m word corpus of English-French sentence pairs or a 91m word corpus of English-German sentence pairs.<sup>9</sup>

To explore the properties of bilingual embeddings learned via objectives other than direct translation, we trained the *BiCVM* model of [50] on the same data, and also downloaded the projected embeddings of [37], *FD*, trained on a bilingual corpus of comparable size ( $\approx 300$  million words per language).<sup>10</sup> Finally, for an initial comparison with monolingual models, we trained a conventional skipgram model [84] and its *Glove* variant [97] for the same number of epochs on the English half of the bilingual corpus.

To analyse the effect on embedding quality of increasing the quantity of training data, we then trained the monolingual models on increasingly large random subsamples of Wikipedia text (up to a total of 1.1bn words). Lastly, we extracted embeddings from a full-sentence language model, *CW*, [28], which was trained for several months on the same Wikipedia 1bn word corpus. Note that increasing the volume of training data for the bilingual (and NMT) models was not possible because of the limited size of available sentence-aligned bitexts.

---

<sup>8</sup>Access to source code and limited GPU time prevent us from training and evaluating the embeddings from other NMT models such as that of [63], [33] and [110]. The underlying principles of encoding-decoding also apply to these models, and we expect the embeddings would exhibit similar properties to those analysed here.

<sup>9</sup>These corpora were produced from the WMT '14 parallel data after conducting the data-selection procedure described by [26].

<sup>10</sup>Available from <http://www.cs.cmu.edu/~mfaruqui/soft.html>. The available embeddings were trained on English-German aligned data, but the authors report similar to for English-French.

### 3.7.1 Similarity and relatedness modelling

As in previous studies [1, 20, 7], our initial evaluations involved calculating pairwise (cosine) distances between embeddings and correlating these distances with (gold-standard) human judgements of the strength of relationships between concepts. For this we used three different gold standards: WordSim-353 [1], MEN [20] and SimLex-999 [54]. Importantly, there is a clear distinction between WordSim-353 and MEN, on the one hand, and SimLex-999, on the other, in terms of the semantic relationship that they quantify. For both WordSim-353 and MEN, annotators were asked to rate how *related* or *associated* two concepts are. Consequently, pairs such as [*clothes-closet*], which are clearly related but ontologically dissimilar, have high ratings in WordSim-353 and MEN. In contrast, such pairs receive a low rating in SimLex-999, where only genuinely *similar* concepts, such as [*coast-shore*], receive high ratings.

To reproduce the scores in SimLex-999, models must thus distinguish pairs that are similar from those that are merely related. In particular, this requires models to develop sensitivity to the distinction between synonyms (similar) and antonyms (often strongly related, but highly dissimilar).<sup>11</sup>

Table 3.4 shows the correlations of NMT (English-French) embeddings, other bilingually-trained embeddings and monolingual embeddings with these three lexical gold-standards. NMT outperform monolingual embeddings, and, to a lesser extent, the other bilingually trained embeddings, on SimLex-999. However, this clear advantage is not observed on MEN and WordSim-353, where the projected embeddings of [37], which were tuned for high performance on WordSim-353, perform best. Given the aforementioned differences between the evaluations, this suggests that bilingually-trained embeddings, and NMT based embeddings in particular, better capture similarity, whereas monolingual embedding spaces are orientated more towards relatedness.

To test this hypothesis further, we ran three more evaluations designed to probe the sensitivity of models to similarity as distinct from relatedness or association. In the first, we measured performance on SimLex-Assoc-333 [54]. This evaluation comprises the 333 most related pairs in SimLex-999, according to an independent empirical measure of relatedness (free associate generation [92]). Importantly, the pairs in SimLex-Assoc-333, while all strongly related, still span the full range of similarity scores.<sup>12</sup> Therefore, the extent to which embeddings can model this data reflects their

---

<sup>11</sup>For a more detailed discussion of the similarity/relatedness distinction, see [54].

<sup>12</sup>The most dissimilar pair in SimLex-Assoc-333 is [*shrink, grow*] with a score of 0.23. The highest is [*vanish, disappear*] with 9.80.

		Monolingual models			Biling. models		NMT models	
		<b>Skipgram</b>	<b>Glove</b>	<b>CW</b>	<b>FD</b>	<b>BiCVM</b>	<b>RNNenc</b>	<b>RNNsearch</b>
WordSim-353	$\rho$	0.52	0.55	0.51	<b>0.69</b>	0.50	0.57	0.58
MEN	$\rho$	0.44	0.71	0.60	<b>0.78</b>	0.45	0.63	0.62
SimLex-999	$\rho$	0.29	0.32	0.28	0.39	0.36	<b>0.52</b>	0.49
SimLex-333	$\rho$	0.18	0.18	0.07	0.24	0.34	<b>0.49</b>	0.45
TOEFL	%	0.75	0.78	0.64	0.84	0.87	<b>0.93</b>	<b>0.93</b>
Syn/antonym	%	0.69	0.72	0.75	0.76	0.70	<b>0.79</b>	0.74

**Table 3.4:** NMT embeddings (RNNenc and RNNsearch) clearly outperform alternative embedding-learning architectures on tasks that require modelling similarity (below the dashed line), but not on tasks that reflect relatedness. Bilingual embedding spaces learned without the translation objective are somewhere between these two extremes.

	<b>Skipgram</b>	<b>Glove</b>	<b>CW</b>	<b>FD</b>	<b>BiCVM</b>	<b>RNNenc</b>	<b>RNNsearch</b>
<i>teacher</i>	<i>vocational</i>	<i>student</i>	<i>student</i>	<i>elementary</i>	<i>faculty</i>	<i>professor</i>	<i>instructor</i>
	<i>in-service</i>	<i>pupil</i>	<i>tutor</i>	<i>school</i>	<i>professors</i>	<i>instructor</i>	<i>professor</i>
	<i>college</i>	<i>university</i>	<i>mentor</i>	<i>classroom</i>	<i>teach</i>	<i>trainer</i>	<i>educator</i>
<i>eaten</i>	<i>spoiled</i>	<i>cooked</i>	<i>baked</i>	<i>ate</i>	<i>eating</i>	<i>ate</i>	<i>ate</i>
	<i>squeezed</i>	<i>eat</i>	<i>peeled</i>	<i>meal</i>	<i>eat</i>	<i>consumed</i>	<i>consumed</i>
	<i>cooked</i>	<i>eating</i>	<i>cooked</i>	<i>salads</i>	<i>baking</i>	<i>tasted</i>	<i>eat</i>
<i>Britain</i>	<i>Northern</i>	<i>Ireland</i>	<i>Luxembourg</i>	<i>UK</i>	<i>UK</i>	<i>UK</i>	<i>England</i>
	<i>Great</i>	<i>Kingdom</i>	<i>Belgium</i>	<i>British</i>	<i>British</i>	<i>British</i>	<i>UK</i>
	<i>Ireland</i>	<i>Great</i>	<i>Madrid</i>	<i>London</i>	<i>England</i>	<i>America</i>	<i>Syria</i>

**Table 3.5:** Nearest neighbours (excluding plurals) in the embedding spaces of different models. All models were trained for 6 epochs on the translation corpus except CW and FD (as noted previously). NMT embedding spaces are oriented according to similarity, whereas embeddings learned by monolingual models are organized according to relatedness. The other bilingual model BiCVM also exhibits a notable focus on similarity.

sensitivity to the similarity (or dissimilarity) of two concepts, even in the face of a strong signal in the training data that those concepts are related.

The TOEFL synonym test is another similarity-focused evaluation of embedding spaces. This test contains 80 cue words, each with four possible answers, of which one is a correct synonym [69]. We computed the proportion of questions answered correctly by each model, where a model’s answer was the nearest (cosine) neighbour to the cue word in its vocabulary.<sup>13</sup> Note that, since TOEFL is a test of synonym recognition, it necessarily requires models to recognise similarity as opposed to relatedness.

Finally, we tested how well different embeddings enabled a supervised classifier

<sup>13</sup>To control for different vocabularies, we restricted the effective vocabulary of each model to the intersection of all model vocabularies, and excluded all questions that contained an answer outside of this intersection.



to distinguish between synonyms and antonyms, since synonyms are necessarily similar and people often find antonyms, which are necessarily dissimilar, to be strongly associated. For 744 word pairs hand-selected as either synonyms or antonyms,<sup>14</sup> we presented a Gaussian SVM with the concatenation of the two word embeddings. We evaluated accuracy using 10-fold cross-validation.

As shown in Table 3.4, with these three additional similarity-focused tasks we again see the same pattern of results. NMT embeddings outperform other bilingually-trained embeddings which in turn outperform monolingual models. The difference is particularly striking on SimLex-Assoc-333, which suggests that the ability to discern similarity from relatedness (when relatedness is high) is perhaps the most clear distinction between the bilingual spaces and those of monolingual models.

These conclusions are also supported by qualitative analysis of the various embedding spaces. As shown in Table 3.5, in the NMT embedding spaces the nearest neighbours (by cosine distance) to concepts such as *teacher* are genuine synonyms such as *professor* or *instructor*. The bilingual objective also seems to orientate the non-NMT embeddings towards semantic similarity, although some purely related neighbours are also observed. In contrast, in the monolingual embedding spaces the neighbours of *teacher* include highly related but dissimilar concepts such as *student* or *college*.

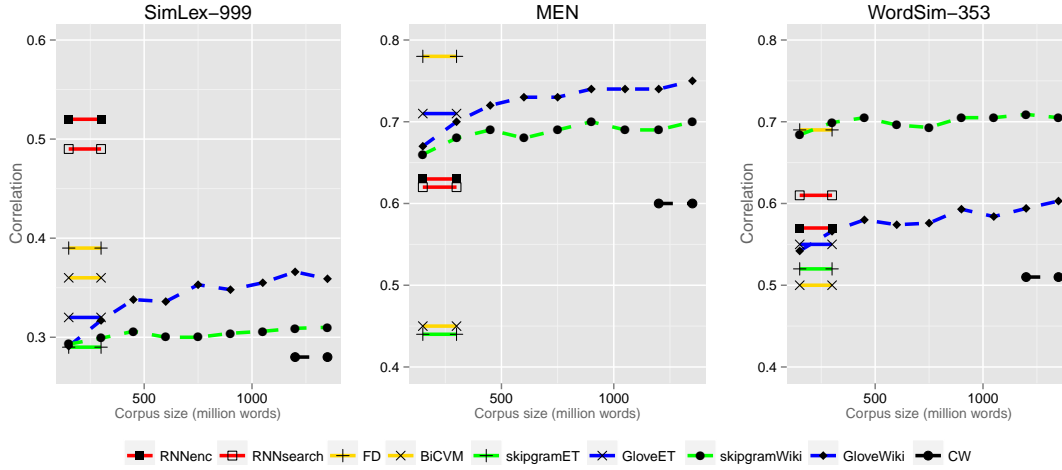
### 3.7.2 Importance of training data quantity

In previous work, monolingual models were trained on corpora many times larger than the English half of our parallel translation corpus. Indeed, the ability to scale to large quantities of training data was one of the principal motivations behind the skipgram architecture [84]. To check if monolingual models simply need more training data to capture similarity as effectively as bilingual models, we therefore trained them on increasingly large subsets of Wikipedia.<sup>15</sup> As shown in Figure 3.5, this is not in fact the case. The performance of monolingual embeddings on similarity tasks remains well below the level of the NMT embeddings and somewhat lower than the non-MT bilingual embeddings as the amount of training data increases.

---

<sup>14</sup>Available online at <http://www.cl.cam.ac.uk/~fh295/>.

<sup>15</sup>We did not do the same for our translation models because sentence-aligned bilingual corpora of comparable size do not exist.



**Figure 3.5:** The effect of increasing the amount of training data on the quality of monolingual embeddings, based on similarity-based evaluations (SimLex-999) and two relatedness-based evaluations (MEN and WordSim-353). *ET* in the legend indicates models trained on the English half of the translation corpus. *Wiki* indicates models trained on Wikipedia.

### 3.7.3 Analogy Resolution

Lexical analogy questions have been used as an alternative way of evaluating word representations. In this task, models must identify the correct answer (*girl*) when presented with analogy questions such as ‘*man* is to *boy* as *woman* is to ?’. It has been shown that Skipgram-style models are surprisingly effective at answering such questions [84]. This is because, if  $\mathbf{m}$ ,  $\mathbf{b}$  and  $\mathbf{w}$  are skipgram-style embeddings for *man*, *boy* and *woman* respectively, the correct answer is often the nearest neighbour in the vocabulary (by cosine distance) to the vector  $\mathbf{v} = \mathbf{w} + \mathbf{b} - \mathbf{m}$ .

We evaluated embeddings on analogy questions using the same vector-algebra method as [84]. As in the previous section, for fair comparison we excluded questions containing a word outside the intersection of all model vocabularies, and restricted all answer searches to this reduced vocabulary. This left 11,166 analogies. Of these, 7219 are classed as ‘syntactic’, in that they exemplify mappings between parts-of-speech or syntactic roles (e.g. *fast* is to *fastest* as *heavy* is to *heaviest*), and 3947 are classed as ‘semantic’ (*Ottawa* is to *Canada* as *Paris* is to *France*), since successful answering seems to rely on some (world) knowledge of the concepts themselves.

As shown in Fig. 3.6, NMT embeddings yield relatively poor answers to semantic analogy questions compared with monolingual embeddings and the bilingual embeddings *FD* (which are projections of similar monolingual embeddings).<sup>16</sup> It appears

<sup>16</sup>The performance of the *FD* embeddings on this task is higher than that reported by [37] because we

that the translation objective prevents the embedding space from developing the same linear, geometric regularities as skipgram-style models with respect to semantic organisation. This also seems to be true of the embeddings from the full-sentence language model *CW*. Further, in the case of the Glove and FD models this advantage seems to be independent of both the domain and size of the training data, since embeddings from these models trained on only the English half of the translation corpus still outperform the translation embeddings.

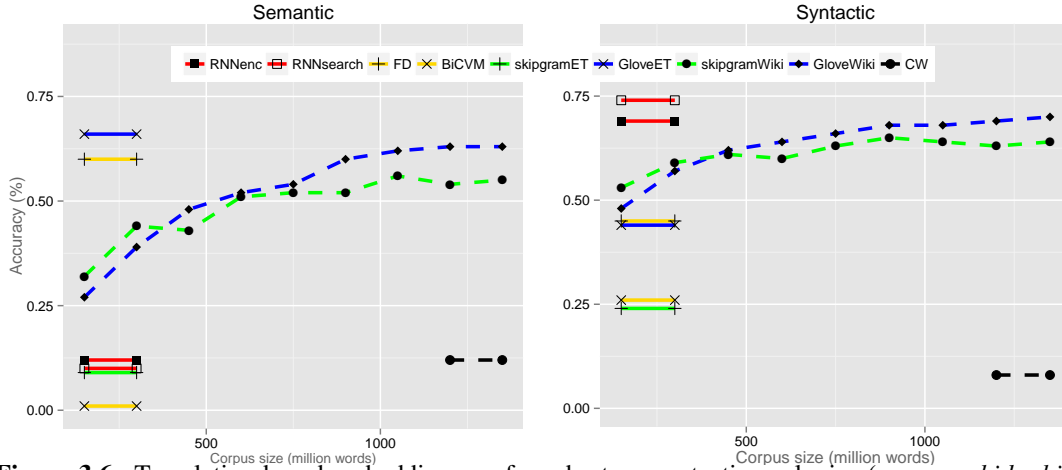
On the other hand, NMT embeddings are effective for answering syntactic analogies using the vector algebra method. They perform comparably to or even better than monolingual embeddings when trained on less data (albeit bilingual data). It is perhaps unsurprising that the translation objective incentivises the encoding of a high degree of lexical syntactic information, since coherent target-language sentences could not be generated without knowledge of the parts-of-speech, tense or case of its vocabulary items. The connection between the translation objective and the embedding of lexical syntactic information is further supported by the fact that embeddings learned by the bilingual model BiCVM do not perform comparably on the syntactic analogy task. In this model, sentential semantics is transferred via a bag-of-words representation, presumably rendering the precise syntactic information less important.

When considering the two properties of NMT embeddings highlighted by these experiments, namely the encoding of semantic similarity and lexical syntax, it is worth noting that items in the similarity-focused evaluations of the previous section (SimLex-999 and TOEFL) consist of word groups or pairs that have identical syntactic role. Thus, even though lexical semantic information is in general pertinent to conceptual similarity [72], the lexical syntactic and conceptual properties of translation embeddings are in some sense independent of one another.

### 3.8 Effect of Target Language

To better understand why a translation objective yields embedding spaces with particular properties, we trained the RNN Search architecture to translate from English to German.

As shown in Table 3.6 (left side), the performance of the source (English) embeddings learned by this model was comparable to that of those learned by the English-to-search for answers over a smaller total candidate vocabulary.



**Figure 3.6:** Translation-based embeddings perform best on syntactic analogies (*run,ran: hide, hid*). Monolingual skipgram/Glove models are better at semantic analogies (*father, man; mother, woman*)

		EN-FR	EN-DE		'earned'	'castle'	'money'
WordSim-353	$\rho$	0.60	<b>0.61</b>	EN-FR	<i>gained</i>	<i>chateau</i>	<b><i>silver</i></b>
MEN	$\rho$	0.61	<b>0.62</b>		<b><i>won</i></b>	<i>palace</i>	<i>funds</i>
SimLex-999	$\rho$	0.49	<b>0.50</b>		<i>acquired</i>	<i>fortress</i>	<i>cash</i>
SimLex-Assoc-333	$\rho$	0.45	<b>0.47</b>	EN-DE	<i>gained</i>	<i>chateau</i>	<i>funds</i>
TOEFL	%	0.90	<b>0.93</b>		<i>deserved</i>	<i>palace</i>	<i>cash</i>
Syn/antonym	%	<b>0.72</b>	0.70		<i>accumulated</i>	<b><i>padlock</i></b>	<i>resources</i>
Syntactic analogies	%	<b>0.73</b>	0.62				
Semantic analogies	%	0.10	<b>0.11</b>				

**Table 3.6:** Comparison of embeddings learned by RNN Search models translating between English-French (EN-FR) and English-German (EN-DE) on all semantic evaluations (left) and nearest neighbours of selected cue words (right). Bold italics indicate target-language-specific effects. Evaluation items and vocabulary searches were restricted to words common to both models.

French model on all evaluations, even though the English-German training corpus (91 million words) was notably smaller than the English-French corpus (348m words). This evidence shows that the desirable properties of translation embeddings highlighted thus far are not particular to English-French translation, and can also emerge when translating to a different language family, with different word ordering conventions.

### 3.9 Overcoming the Vocabulary Size Problem

A potential drawback to using NMT models for learning word embeddings is the computational cost of training such a model on large vocabularies. To generate a target language sentence, NMT models repeatedly compute a softmax distribution over the target vocabulary. This computation scales with vocabulary size and must be repeated for each word in the output sentence, so that training models with large output vocabularies is challenging. Moreover, while the same computational bottleneck does not apply to the encoding process or source vocabulary, there is no way in which a translation model could learn a high quality source embedding for a word if the plausible translations were outside its vocabulary. Thus, limitations on the size of the target vocabulary effectively limit the scope of NMT models as representation-learning tools. This contrasts with the shallower monolingual and bilingual representation-learning models considered in this paper, which efficiently compute a distribution over a large target vocabulary using either a hierarchical softmax [90] or approximate methods such as negative sampling [84, 50], and thus can learn large vocabularies of both source and target embeddings.

A recently proposed solution to this problem enables NMT models to be trained with larger target vocabularies (and hence larger meaningful source vocabularies) at comparable computational cost to training with a small target vocabulary [60]. The algorithm uses (biased) importance sampling [11] to approximate the probability distribution of words over a large target vocabulary with a finite set of distributions over subsets of that vocabulary. Despite this element of approximation in the decoder, extending the effective target vocabulary in this way significantly improves translation performance, since the model can make sense of more sentences in the training data and encounters fewer unknown words at test time. In terms of representation learning, the method provides a means to scale up the NMT approach to vocabularies as large as those learned by monolingual models. However, given that the method replaces an exact calculation with an approximate one, we tested how the quality of source embeddings is affected by scaling up the target language vocabulary in this way.

As shown in Table 3.7, there is no significant degradation of embedding quality when scaling to large vocabularies with using the approximate decoder. Note that for a fair comparison we filtered these evaluations to only include items that are present in the smaller vocabulary. Thus, the numbers do not directly reflect the quality of the additional 470k embeddings learned by the extended vocabulary models, which

		RNN Search	RNN Search	RNN Search-LV	RNN Search-LV
		EN-FR	EN-DE	EN-FR	EN-DE
WordSim-353	$\rho$	0.60	<b>0.61</b>	0.59	0.57
MEN	$\rho$	0.61	<b>0.62</b>	<b>0.62</b>	0.61
SimLex-999	$\rho$	0.49	0.50	<b>0.51</b>	0.50
SimLex-Assoc-333	$\rho$	0.45	<b>0.47</b>	<b>0.47</b>	0.46
TOEFL	%	0.90	0.93	0.93	<b>0.98</b>
Syn/antonym	%	0.72	0.70	<b>0.74</b>	0.71
Syntactic analogies	%	<b>0.73</b>	0.62	0.71	0.62
Semantic analogies	%	0.10	0.11	0.08	<b>0.13</b>

**Table 3.7:** Comparison of embeddings learned by the original (RNN Search - 30k French words, 50k German words) and extended-vocabulary (RNN Search-LV -500k words) models translating from English to French (EN-FR) and from English to German (EN-DE). For fair comparisons, all evaluations were restricted to the intersection of all model vocabularies.

one would expect to be lower since they are words of lower frequency. All embeddings can be downloaded from <http://www.cl.cam.ac.uk/~fh295/>, and the embeddings from the smaller vocabulary models can be interrogated at <http://lisa.iro.umontreal.ca/mt-demo/embs/>.<sup>17</sup>

### 3.10 How Similarity Emerges

Although NMT models appear to encode both conceptual similarity and syntactic information for any source and target languages, it is not the case that embedding spaces will always be identical. Interrogating the nearest neighbours of the source embedding spaces of the English-French and English-German models reveals occasional language-specific effects. As shown in Table 3.6 (right side), the neighbours for the word *earned* in the English-German model are as one might expect, whereas the neighbours from the English-French model contain the somewhat unlikely candidate *won*. In a similar vein, while the neighbours of the word *castle* from the English-French model are unarguably similar, the neighbours from the English-German model contain the word *padlock*.

These infrequent but striking differences between the English-German and English-French source embedding spaces indicate how similarity might emerge effectively in NMT models. Tokens of the French verb *gagner* have (at least) two possible English

<sup>17</sup>A different solution to the rare-word problem was proposed by [76]. We do not evaluate the effects on the resulting embeddings of this method because we lack access to the source code.

translations (*win* and *earn*). Since the translation model, which has limited encoding capacity, is trained to map tokens of *win* and *earn* to the same place in the target embedding space, it is efficient to move these concepts closer in the source space. Since *win* and *earn* map directly to two different verbs in German, this effect is not observed. On the other hand, the English nouns *castle* and *padlock* translate to a single noun (*Schloss*) in German, but different nouns in French. Thus, *padlock* and *castle* are only close in the source embeddings from the English-German model.

Based on these considerations, we can conjecture that the following condition on the semantic configuration between two language is crucial to the effective induction of lexical similarity.

- (1) For  $s_1$  and  $s_2$  in the source language, there is some  $t$  in the target language such that there are sentences in the training data in which  $s_1$  translates to  $t$  and sentences in which  $s_2$  translates to  $t$ .

*if and only if*

- (2)  $s_1$  and  $s_2$  are semantically similar.

Of course, this condition is not true in general. However, we propose that the extent to which it holds over all possible word pairs corresponds to the quality of similarity induction in the translation embedding space. Note that strong polysemy in the target language, such as *gagner* = *win*, *earn*, can lead to cases in which 1 is satisfied but 2 is not. The conjecture claims that these cases are detrimental to the quality of the embedding space (at least with regards to similarity). In practice, qualitative analyses of the embedding spaces and native speaker intuitions suggest that such cases are comparatively rare. Moreover, when such cases are observed,  $s_1$  and  $s_2$ , while perhaps not similar, are not strongly dissimilar. This could explain why related but strongly dissimilar concepts such as antonym pairs do not converge in the translation embedding space. This is also consistent with qualitative evidence presented by [37] that projecting monolingual embeddings into a bilingual space orientates them to better reflect the synonymy/antonymy distinction.

## 3.11 Conclusion

In this work, we have shown that the embedding spaces from neural machine translation models are orientated more towards conceptual similarity than those of monolingual models, and that translation embedding spaces also reflect richer lexical syntactic information. To perform well on similarity evaluations such as SimLex-999, embeddings must distinguish information pertinent to what concepts *are* (their function or ontology) from information reflecting other non-specific inter-concept relationships. Concepts that are strongly related but dissimilar, such as antonyms, are particularly challenging in this regard [54]. Consistent with the qualitative observation made by [37], we suggested how the nature of the semantic correspondence between the words in languages enables NMT embeddings to distinguish synonyms and antonyms and, more generally, to encode the information needed to reflect human intuitions of similarity.

The language-specific effects we observed in Section 3.8 suggest a potential avenue for improving translation and multi-lingual embeddings in future work. First, as the availability of fast GPUs for training grows, we would like to explore the embeddings learned by NMT models that translate between much more distant language pairs such as English-Chinese or English-Arabic. For these language pairs, the word alignment will less monotonic and may result in even more important semantic and syntactic information being encoded in the lexical representation. Further, as observed by both [50] and [37], the bilingual representation learning paradigm can be naturally extended to update representations based on correspondences between multiple languages (for instance by interleaving English-French and English-German training examples). Such an approach should smooth out language-specific effects, leaving embeddings that encode only language-agnostic conceptual semantics and are thus more generally applicable. Another related challenge is to develop smaller or less complex representation-learning tools that encode similarity with as much fidelity as NMT models but without the computational overhead. One promising approach for this is to learn word alignments and word embeddings jointly [68]. This approach is effective for cross-lingual document classification, although the authors do evaluate the monolingual subspace induced by the model.<sup>18</sup>

Not all word embeddings learned from text are born equal. Depending on the ap-

---

<sup>18</sup>These embeddings are not publicly available and we were unable to re-train them using the source code.



plication, those learned by NMT models may have particularly desirable properties. For decades, distributional semantic models have aimed to exploit Firth’s famous *distributional hypothesis* to induce word meanings from (monolingual) text. However, the hypothesis also betrays the weakness of the monolingual distributional approach when it comes to learning human-quality concept representations. For while it is undeniable that “words which are similar in meaning appear in similar distributional contexts” [40], the converse assertion, which is what really matters, is only sometimes true.



# Chapter 4

## Learning to Represent Phrases

### 4.1 Introduction

Much recent research in computational semantics has focussed on learning representations of arbitrary-length phrases and sentences. This task is challenging partly because there is no obvious gold standard of phrasal representation that could be used in training and evaluation. Consequently, it is difficult to design approaches that could learn from such a gold standard, and also hard to evaluate or compare different models.

In this work, we use dictionary definitions to address this issue. The composed meaning of the words in a dictionary definition (*a tall, long-necked, spotted ruminant of Africa*) should correspond to the meaning of the word they define (*giraffe*). This bridge between lexical and phrasal semantics is useful because high quality vector representations of single words can be used as a target when learning to combine the words into a coherent phrasal representation.

This approach still requires a model capable of learning to map between arbitrary-length phrases and fixed-length continuous-valued word vectors. For this purpose we experiment with two broad classes of neural language models (NLMs): Recurrent Neural Networks (RNNs), which naturally encode the order of input words, and simpler (feedforward) bag-of-words (BOW) embedding models. Prior to training these NLMs, we learn target lexical representations by training the Word2Vec software [84] on billions of words of raw text.

We demonstrate the usefulness of our approach by building and releasing two applications. The first is a *reverse dictionary* or *concept finder*: a system that returns words based on user descriptions or definitions [125]. Reverse dictionaries are used

by copywriters, novelists, translators and other professional writers to find words for notions or ideas that might be on the tip of their tongue. For instance, a travel-writer might look to enhance her prose by searching for examples of a *country that people associate with warm weather* or *an activity that is mentally or physically demanding*. We show that an NLM-based reverse dictionary trained on only a handful of dictionaries identifies novel definitions and concept descriptions comparably or better than commercial systems, which rely on significant task-specific engineering and access to much more dictionary data. Moreover, by exploiting models that learn bilingual word representations [119, 67, 49, 43], we show that the NLM approach can be easily extended to produce a potentially useful cross-lingual reverse dictionary.

The second application of our models is as a general-knowledge crossword question answerer. When trained on both dictionary definitions and the opening sentences of Wikipedia articles, NLMs produce plausible answers to (non-cryptic) crossword clues, even those that apparently require detailed world knowledge. Both BOW and RNN models can outperform bespoke commercial crossword solvers, particularly when clues contain a greater number of words. Qualitative analysis reveals that NLMs can learn to relate concepts that are not directly connected in the training data and can thus generalise well to unseen input. To facilitate further research, all of our code, training and evaluation sets (together with a system demo) are published online with this paper.<sup>1</sup>

## 4.2 Neural Language Model Architectures

The first model we apply to the dictionary-based learning task is a recurrent neural network (RNN). RNNs operate on variable-length sequences of inputs; in our case, natural language definitions, descriptions or sentences. RNNs (with LSTMs) have achieved state-of-the-art performance in language modelling [82], image caption generation [65] and approach state-of-the-art performance in machine translation [5].

During training, the input to the RNN is a dictionary definition or sentence from an encyclopedia. The objective of the model is to map these defining phrases or sentences to an embedding of the word that the definition defines. The target word embeddings are learned independently of the RNN weights, using the Word2Vec software [84].

The set of all words in the training data constitutes the vocabulary of the RNN.

---

<sup>1</sup> <https://www.cl.cam.ac.uk/~fh295/>

For each word in this vocabulary we randomly initialise a real-valued vector (input embedding) of model parameters. The RNN ‘reads’ the first word in the input by applying a non-linear projection of its embedding  $v_1$  parameterised by input weight matrix  $W$  and  $b$ , a vector of biases.

$$A_1 = \phi(Wv_1 + b)$$

yielding the first internal activation state  $A_1$ . In our implementation, we use  $\phi(x) = \tanh(x)$ , though in theory  $\phi$  can be any differentiable non-linear function. Subsequent internal activations (after time-step  $t$ ) are computed by projecting the embedding of the  $t^{th}$  word and using this information to ‘update’ the internal activation state.

$$A_t = \phi(UA_{t-1} + Wv_t + b).$$

As such, the values of the final internal activation state units  $A_N$  are a weighted function of all input word embeddings, and constitute a ‘summary’ of the information in the sentence.

### 4.2.1 Long Short Term Memory

A known limitation when training RNNs to read language using gradient descent is that the error signal (gradient) on the training examples either vanishes or explodes as the number of time steps (sentence length) increases [12]. Consequently, after reading longer sentences the final internal activation  $A_N$  typically retains useful information about the most recently read (sentence-final) words, but can neglect important information near the start of the input sentence. LSTMs [55] were designed to mitigate this long-term dependency problem.

At each time step  $t$ , in place of the single internal layer of units  $A$ , the LSTM RNN computes six internal layers  $i^w, g^i, g^f, g^o, h$  and  $m$ . The first,  $g^w$ , represents the core information passed to the LSTM unit by the latest input word at  $t$ . It is computed as a simple linear projection of the input embedding  $v_t$  (by input weights  $W_w$ ) and the *output state* of the LSTM at the previous time step  $h_{t-1}$  (by update weights  $U_w$ ):

$$i_t^w = W_w v_t + U_w h_{t-1} + b_w$$

The layers  $g^i, g^f$  and  $g^o$  are computed as weighted sigmoid functions of the input

embeddings, again parameterised by layer-specific weight matrices  $W$  and  $U$ :

$$g_t^s = \frac{1}{1 + \exp(-(W_s v_t + U_s h_{t-1} + b_s))}$$

where  $s$  stands for one of  $i$ ,  $f$  or  $o$ . These vectors take values on  $[0, 1]$  and are often referred to as *gating activations*. Finally, the *internal memory state*,  $m_t$  and new output state  $h_t$ , of the LSTM at  $t$  are computed as

$$\begin{aligned} m_t &= i_t^w \odot g_t^i + m_{t-1} \odot g_t^f \\ h_t &= g_t^o \odot \phi(m_t), \end{aligned}$$

where  $\odot$  indicates elementwise vector multiplication and  $\phi$  is, as before, some non-linear function (we use *tanh*). Thus,  $g^i$  determines to what extent the new *input* word is considered at each time step,  $g^f$  determines to what extent the existing state of the internal memory is retained or *forgotten* in computing the new internal memory, and  $g^o$  determines how much this memory is considered when computing the output state at  $t$ .

The sentence-final memory state of the LSTM,  $m_N$ , a ‘summary’ of all the information in the sentence, is then projected via an extra non-linear projection (parameterised by a further weight matrix) to a target embedding space. This layer enables the target (defined) word embedding space to take a different dimension to the activation layers of the RNN, and in principle enables a more complex definition-reading function to be learned.

### 4.2.2 Bag-of-Words NLMs

We implement a simpler linear bag-of-words (BOW) architecture for encoding the definition phrases. As with the RNN, this architecture learns an embedding  $v_i$  for each word in the model vocabulary, together with a single matrix of input projection weights  $W$ . The BOW model simply maps an input definition with word embeddings  $v_1 \dots v_n$  to the sum of the projected embeddings  $\sum_{i=1}^n W v_i$ . This model can also be considered a special case of an RNN in which the update function  $U$  and nonlinearity  $\phi$  are both the identity, so that ‘reading’ the next word in the input phrase updates the current

representation more simply:

$$A_t = A_{t-1} + Wv_t.$$

### 4.2.3 Pre-trained Input Representations

We experiment with variants of these models in which the input definition embeddings are pre-learned and fixed (rather than randomly-initialised and updated) during training. There are several potential advantages to taking this approach. First, the word embeddings are trained on massive corpora and may therefore introduce additional linguistic or conceptual knowledge to the models. Second, at test time, the models will have a larger effective vocabulary, since the pre-trained word embeddings typically span a larger vocabulary than the union of all dictionary definitions used to train the model. Finally, the models will then map to and from the same space of embeddings (the embedding space will be closed under the operation of the model), so conceivably could be more easily applied as a general-purpose ‘composition engine’.

### 4.2.4 Training Objective

We train all neural language models  $M$  to map the input definition phrase  $s_c$  defining word  $c$  to a location close to the the pre-trained embedding  $v_c$  of  $c$ . We experiment with two different cost functions for the word-phrase pair  $(c, s_c)$  from the training data. The first is simply the cosine distance between  $M(s_c)$  and  $v_c$ . The second is the rank loss

$$\max(0, m - \cos(M(s_c), v_c) - \cos(M(s_c), v_r))$$

where  $v_r$  is the embedding of a randomly-selected word from the vocabulary other than  $c$ . This loss function was used for language models, for example, in [57]. In all experiments we apply a margin  $m = 0.1$ , which has been shown to work well on word-retrieval tasks [18].

### 4.2.5 Implementation Details

Since training on the dictionary data took 6-10 hours, we did not conduct a hyperparameter search on any validation sets over the space of possible model configurations such as embedding dimension, or size of hidden layers. Instead, we chose these

parameters to be as standard as possible based on previous research. For fair comparison, any aspects of model design that are not specific to a particular class of model were kept constant across experiments.

The pre-trained word embeddings used in all of our models (either as input or target) were learned by a continuous bag-of-words (CBOW) model using the Word2Vec software on approximately 8 billion words of running text.<sup>2</sup> When training such models on massive corpora, a large embedding length of up to 700 have been shown to yield best performance (see e.g. [36]). The pre-trained embeddings used in our models were of length 500, as a compromise between quality and memory constraints.

In cases where the word embeddings are learned during training on the dictionary objective, we make these embeddings shorter (256), since they must be learned from much less language data. In the RNN models, and at each time step each of the four LSTM RNN internal layers (gating and activation states) had length 512 – another standard choice (see e.g. [25]). The final hidden state was mapped linearly to length 500, the dimension of the target embedding. In the BOW models, the projection matrix projects input embeddings (either learned, of length 256, or pre-trained, of length 500) to length 500 for summing.

All models were implemented with Theano [14] and trained with minibatch SGD on GPUs. The batch size was fixed at 16 and the learning rate was controlled by *adadelta* [124].

### 4.3 Reverse Dictionaries

The most immediate application of our trained models is as a *reverse dictionary* or *concept finder*. It is simple to look up a definition in a dictionary given a word, but professional writers often also require suitable words for a given idea, concept or definition.<sup>3</sup> Reverse dictionaries satisfy this need by returning candidate words given a phrase, description or definition. For instance, when queried with the phrase *an activity that requires strength and determination*, the OneLook.com reverse dictionary returns the concepts *exercise* and *work*. Our trained RNN model can perform a similar function, simply by mapping a phrase to a point in the target (Word2Vec) embedding

---

<sup>2</sup>The Word2Vec embedding models are well known; further details can be found at <https://code.google.com/p/word2vec/>. The training data for this pre-training was compiled from various online text sources using the script *demo-train-big-model-v1.sh* from the same page.

<sup>3</sup>See the testimony from professional writers at <http://www.onelook.com/?c=awards>



space, and returning the words corresponding to the embeddings that are closest to that point.

Several other academic studies have proposed reverse dictionary models. These generally rely on common techniques from information retrieval, comparing definitions in their internal database to the input query, and returning the word whose definition is ‘closest’ to that query [15, 16, 125]. Proximity is quantified differently in each case, but is generally a function of hand-engineered features of the two sentences. For instance, [103] propose a method in which the candidates for a given input query are all words in the model’s database whose definitions contain one or more words from the query. This candidate list is then ranked according to a query-definition similarity metric based on the hypernym and hyponym relations in WordNet, features commonly used in IR such as *tf-idf* and a parser.

There are, in addition, at least two commercial online reverse dictionary applications, whose architecture is proprietary knowledge. The first is the Dictionary.com reverse dictionary <sup>4</sup>, which retrieves candidate words from the Dictionary.com dictionary based on user definitions or descriptions. The second is **OneLook.com**, whose algorithm searches 1061 indexed dictionaries, including all major freely-available online dictionaries and resources such as Wikipedia and WordNet.

### 4.3.1 Data Collection and Training

To compile a bank of dictionary definitions for training the model, we started with all words in the target embedding space. For each of these words, we extracted dictionary-style definitions from five electronic resources: *Wordnet*, *The American Heritage Dictionary*, *The Collaborative International Dictionary of English*, *Wiktionary* and *Webster’s*. We chose these five dictionaries because they are freely-available via the WordNik API,<sup>5</sup> but in theory any dictionary could be chosen. Most words in our training data had multiple definitions. For each word  $w$  with definitions  $\{d_1 \dots d_n\}$  we included all pairs  $(w, d_1) \dots (w, d_n)$  as training examples.

To allow models access to more factual knowledge than might be present in a dictionary (for instance, information about specific entities, places or people, we supplemented this training data with information extracted from Simple Wikipedia.<sup>6</sup> For

---

<sup>4</sup>Available at <http://dictionary.reference.com/reverse/>

<sup>5</sup>See <http://developer.wordnik.com>

<sup>6</sup>[https://simple.wikipedia.org/wiki/Main\\_Page](https://simple.wikipedia.org/wiki/Main_Page)

every word in the model’s target embedding space that is also the title of a Wikipedia article, we treat the sentences in the first paragraph of the article as if they were (independent) definitions of that word. When a word in Wikipedia also occurs in one (or more) of the five training dictionaries, we simply add these pseudo-definitions to the training set of definitions for the word. Combining Wikipedia and dictionaries in this way resulted in  $\approx 900,000$  word-’definition’ pairs of  $\approx 100,000$  unique words.

To explore the effect of the quantity of training data on the performance of the models, we also trained models on subsets of this data. The first subset comprised only definitions from Wordnet (approximately 150,000 definitions of 75,000 words). The second subset comprised only words in Wordnet and their *first* definitions (approximately 75,000 word, definition pairs).<sup>7</sup> For all variants of RNN and BOW models, however, reducing the training data in this way resulted in a clear reduction in performance on all tasks. For brevity, we therefore do not present these results in what follows.

Test Set		Dictionary definitions						Concept descriptions (200)		
		Seen (500 WN defs)			Unseen (500 WN defs)					
Unsup. models	W2V add	-	-	-	923	.04/.16	163	339	.07/.30	150
	W2V mult	-	-	-	1000	.00/.00	10*	1000	.00/.00	27*
	OneLook	<b>0</b>	<b>.89/.91</b>	<b>67</b>	-	-	-	<b>18.5</b>	<b>.38/.58</b>	153
NLMs	RNN cosine	12	.48/.73	103	22	.41/.70	116	69	.28/.54	157
	RNN w2v cosine	19	.44/.70	111	19	.44/.69	126	26	<b>.38/.66</b>	111
	RNN ranking	18	.45/.67	128	24	.43/.69	103	25	.34/.66	102
	RNN w2v ranking	54	.32/.56	155	33	.36/.65	137	30	.33/.69	<b>77</b>
	BOW cosine	22	.44/.65	129	19	.43/.69	103	50	.34/.60	99
	BOW w2v cosine	15	.46/.71	124	<b>14</b>	<b>.46/.71</b>	104	28	.36/.66	99
	BOW ranking	17	.45/.68	115	22	.42/.70	<b>95</b>	32	.35/.69	101
	BOW w2v rankng	55	.32/.56	155	36	.35/.66	138	38	<b>.33/.72</b>	85

<

| median rank    accuracy@10/100    rank variance |

**Table 4.1:** Performance of different reverse dictionary models in different evaluation settings. \*Low variance in *mult* models is due to consistently poor scores, so not highlighted.

### 4.3.2 Comparisons

As a baseline, we also implemented two entirely unsupervised methods using the neural (Word2Vec) word embeddings from the target word space. In the first (**W2V add**),

<sup>7</sup>As with other dictionaries, the first definition in WordNet generally corresponds to the most typical or common sense of a word.

we compose the embeddings for each word in the input query by pointwise addition, and return as candidates the nearest word embeddings to the resulting composed vector.<sup>8</sup> The second baseline, (**W2V mult**), is identical except that the embeddings are composed by elementwise multiplication. Both methods are established ways of building phrase representations from word embeddings [87].

None of the models or evaluations from previous academic research on reverse dictionaries is publicly available, so direct comparison is not possible. However, we do compare performance with the commercial systems. The Dictionary.com system returned no candidates for over 96% of our input definitions. We therefore conduct detailed comparison with OneLook.com, which is the first reverse dictionary tool returned by a Google search and seems to be the most popular among writers.

### 4.3.3 Reverse Dictionary Evaluation

To our knowledge there are no established means of measuring reverse dictionary performance. In the only previous academic research on English reverse dictionaries that we are aware of, evaluation was conducted on 300 word-definition pairs written by lexicographers [103]. Since these are not publicly available we developed new evaluation sets and make them freely available for future evaluations.

The evaluation items are of three types, designed to test different properties of the models. To create the **seen** evaluation, we randomly selected 500 words from the WordNet training data (seen by all models), and then randomly selected a definition for each word. Testing models on the resulting 500 word-definition pairs assesses their ability to recall or decode previously encoded information. For the **unseen** evaluation, we randomly selected 500 words from WordNet and excluded all definitions of these words from the training data of all models.

Finally, for a fair comparison with OneLook, which has both the seen and unseen pairs in its internal database, we built a new dataset of **concept descriptions** that do not appear in the training data for any model. To do so, we randomly selected 200 adjectives, nouns or verbs from among the top 3000 most frequent tokens in the British National Corpus [71] (but outside the top 100). We then asked ten native English speakers to write a single-sentence ‘description’ of these words. To ensure the resulting descriptions were good quality, for each description we asked two participants

---

<sup>8</sup>Since we retrieve all answers from embedding spaces by cosine similarity, addition of word embeddings is equivalent to taking the mean.

who did not produce that description to list any words that fitted the description (up to a maximum of three). If the target word was not produced by one of the two checkers, the original participant was asked to re-write the description until the validation was passed.<sup>9</sup> These concept descriptions, together with other evaluation sets, can be downloaded from our website for future comparisons.

Test set	Word	Description
Dictionary definition	<i>valve</i>	"control consisting of a mechanical device for controlling fluid flow"
Concept description	<i>prefer</i>	"when you like one thing more than another thing"

**Table 4.2:** Style difference between *dictionary definitions* and *concept descriptions* in the evaluation.

Given a test description, definition, or question, all models produce a ranking of possible word answers based on the proximity of their representations of the input phrase and all possible output words. To quantify the quality of a given ranking, we report three statistics: the *median rank* of the correct answer (over the whole test set, lower better), the proportion of training cases in which the correct answer appears in the top 10/100 in this ranking (*accuracy@10/100* - higher better) and the variance of the rank of the correct answer across the test set (*rank variance* - lower better).

#### 4.3.4 Results

Table 4.1 shows the performance of the different models in the three evaluation settings. Of the unsupervised composition models, elementwise addition is clearly more effective than multiplication, which almost never returns the correct word as the nearest neighbour of the composition. Overall, however, the supervised models (RNN, BOW and OneLook) clearly outperform these baselines.

The results indicate interesting differences between the NLMs and the OneLook dictionary search engine. The Seen (WN first) definitions in Table 4.1 occur in both the training data for the NLMs and the lookup data for the OneLook model. Clearly the OneLook algorithm is better than NLMs at retrieving already available information (returning 89% of correct words among the top-ten candidates on this set). However, this is likely to come at the cost of a greater memory footprint, since the model requires access to its database of dictionaries at query time.<sup>10</sup>

<sup>9</sup>Re-writing was required in 6 of the 200 cases.

<sup>10</sup>The trained neural language models are approximately half the size of the six training dictionary-

The performance of the NLM embedding models on the (unseen) concept descriptions task shows that these models can generalise well to novel, unseen queries. While the median rank for OneLook on this evaluation is lower, the NLMs retrieve the correct answer in the top ten candidates approximately as frequently, within the top 100 candidates more frequently and with lower variance in ranking over the test set. Thus, NLMs seem to generalise more ‘consistently’ than OneLook on this dataset, in that they generally assign a reasonably high ranking to the correct word. In contrast, as can also be verified by querying our web demo, OneLook tends to perform either very well or poorly on a given query.<sup>11</sup>

When comparing between NLMs, perhaps the most striking observation is that the RNN models do not significantly outperform the BOW models, even though the BOW model output is invariant to changes in the order of words in the definition. Users of the online demo can verify that the BOW models recover concepts from descriptions strikingly well, even when the words in the description are permuted. This observation underlines the importance of lexical semantics in the interpretation of language by NLMs, and is consistent with some other recent work on embedding sentences [59].

It is difficult to observe clear trends in the differences between NLMs that learn input word embeddings and those with pre-trained (Word2Vec) input embeddings. Both types of input yield good performance in some situations and weaker performance in others. In general, pre-training input embeddings seems to help most on the concept descriptions, which are furthest from the training data in terms of linguistic style. This is perhaps unsurprising, since models that learn input embeddings from the dictionary data acquire all of their conceptual knowledge from this data (and thus may overfit to this setting), whereas models with pre-trained embeddings have some semantic memory acquired from general running-text language data and other knowledge acquired from the dictionaries.

### 4.3.5 Qualitative Analysis

Some example output from the various models is presented in Table 4.3. The differences illustrated here are also evident from querying the web demo. The first example shows how the NLMs (BOW and RNN) generalise beyond their training data. Four of

---

ies stored as plain text, so would be hundreds of times smaller than the OneLook database of 1061 dictionaries if stored this way.

<sup>11</sup>We also observed that the *mean* ranking for NLMs was lower than for OneLook on the concept descriptions task.

Input Description	OneLook	W2V add	RNN	BOW
"a native of a cold country"	1:country 2:citizen 3:foreign 4:naturalize 5:cisco	1:a 2:the 3:another 4:of 5:whole	1:eskimo 2:scandinavian 3:arctic 4:indian 5:siberian	1:frigid 2:cold 3:icy 4:russian 5:indian
"a way of moving through the air"	1:drag 2:whiz 3:aerodynamics 4:draught 5:coefficient of drag	1:the 2:through 3:a 4:moving 5:in	1:glide 2:scooting 3:glides 4:gliding 5:flight	1:flying 2:gliding 3:glide 4:fly 5:scooting
"a habit that might annoy your spouse"	1:sisterinlaw 2:fatherinlaw 3:motherinlaw 4:stepson 5:stepchild	1:annoy 2:your 3:might 4:that 5:either	1:bossiness 2:jealousy 3:annoyance 4:rudeness 5:boorishness	1:infidelity 2:bossiness 3:foible 4:unfaithfulness 5:adulterous

**Table 4.3:** The top-five candidates for example queries (invented by the authors) from different reverse dictionary models. Both the RNN and BOW models are without Word2Vec input and use the cosine loss.

the top five responses could be classed as appropriate in that they refer to inhabitants of cold countries. However, inspecting the WordNik training data, there is no mention of *cold* or anything to do with climate in the definitions of *Eskimo*, *Scandinavian*, *Scandinavia* etc. Therefore, the embedding models must have learned that *coldness* is a characteristic of Scandinavia, Siberia, Russia, relates to Eskimos etc. via connections with other concepts that are described or defined as *cold*. In contrast, the candidates produced by the OneLook and (unsupervised) W2V baseline models have nothing to do with coldness.

The second example demonstrates how the NLMs generally return candidates whose linguistic or conceptual function is appropriate to the query. For a query referring explicitly to a means, method or process, the RNN and BOW models produce verbs in different forms or an appropriate deverbal noun. In contrast, OneLook returns words of all types (*aerodynamics*, *draught*) that are arbitrarily related to the words in the query. A similar effect is apparent in the third example. While the candidates produced by the OneLook model are the correct part of speech (Noun), and related to the query topic, they are not semantically appropriate. The dictionary embedding models are the only ones that return a list of plausible *habits*, the class of noun requested by the input.

Input description	RNN EN-FR	W2V add	RNN + Google
"an emotion that you might feel after being rejected"	<u>triste</u> , pitoyable <u>répugnante</u> , épouvantable	insister, effectivement pourquoi, nous	sentiment, regret <u>peur</u> , aversion
"a small black flying insect that transmits disease and likes horses"	<u>mouche</u> , canard <u>hirondelle</u> , pigeon	attentivement, pouvions pourrons, naturellement	voler, <u>faucon</u> <u>mouches</u> , volant

**Table 4.4:** Responses from cross-lingual reverse dictionary models to selected queries. Underlined responses are ‘correct’ or potentially useful for a native French speaker.

### 4.3.6 Cross-Lingual Reverse Dictionaries

We now show how the RNN architecture can be easily modified to create a *bilingual reverse dictionary* - a system that returns candidate words in one language given a description or definition in another. A bilingual reverse dictionary could have clear applications for translators or transcribers. Indeed, the problem of attaching appropriate words to concepts may be more common when searching for words in a second language than in a monolingual context.

To create the bilingual variant, we simply replace the Word2Vec target embeddings with those from a bilingual embedding space. Bilingual embedding models use bilingual corpora to learn a space of representations of the words in two languages, such that words from either language that have similar meanings are close together [49, 22, 43]. For a test-of-concept experiment, we used English-French embeddings learned by the state-of-the-art BilBOWA model [43] from the Wikipedia (monolingual) and Europarl (bilingual) corpora.<sup>12</sup> We trained the RNN model to map from English definitions to English words in the bilingual space. At test time, after reading an English definition, we then simply return the nearest French word neighbours to that definition.

Because no benchmarks exist for quantitative evaluation of bilingual reverse dictionaries, we compare this approach qualitatively with two alternative methods for mapping definitions to words across languages. The first is analogous to the W2V Add model of the previous section: in the bilingual embedding space, we first compose the embeddings of the English words in the query definition with elementwise addition, and then return the French word whose embedding is nearest to this vector sum. The second uses the RNN monolingual reverse dictionary model to identify an English word from an English definition, and then translates that word using Google Translate.

Table 4.4 shows that the RNN model can be effectively modified to create a cross-

<sup>12</sup>The approach should work with any bilingual embeddings. We thank Stephan Gouws for doing the training.

lingual reverse dictionary. It is perhaps unsurprising that the W2V Add model candidates are generally the lowest in quality given the performance of the method in the monolingual setting. In comparing the two RNN-based methods, the RNN (embedding space) model appears to have two advantages over the RNN + Google approach. First, it does not require online access to a bilingual word-word mapping as defined e.g. by Google Translate. Second, it is less prone to errors caused by word sense ambiguity. For example, in response to the query *an emotion you feel after being rejected*, the bilingual embedding RNN returns emotions or adjectives describing mental states. In contrast, the monolingual+Google model incorrectly maps the plausible English response *regret* to the verbal infinitive *regretter*. The model makes the same error when responding to a description of a fly, returning the verb *voler* (to fly).

### 4.3.7 Discussion

We have shown that simply training RNN or BOW NLMs on six dictionaries yields a reverse dictionary that performs comparably to the leading commercial system, even with access to much less dictionary data. Indeed, the embedding models consistently return syntactically and semantically plausible responses, which are generally part of a more coherent and homogeneous set of candidates than those produced by the commercial systems. We also showed how the architecture can be easily extended to produce bilingual versions of the same model.

In the analyses performed thus far, we only test the dictionary embedding approach on tasks that it was trained to accomplish (mapping definitions or descriptions to words). In the next section, we explore whether the knowledge learned by dictionary embedding models can be effectively transferred to a novel task.

## 4.4 General Knowledge (crossword) Question Answering

The automatic answering of questions posed in natural language is a central problem of Artificial Intelligence. Although web search and IR techniques provide a means to find sites or documents related to language queries, at present, internet users requiring a specific fact must still sift through pages to locate the desired information.

Systems that attempt to overcome this, via fully open-domain or general knowledge



question-answering (open QA), generally require large teams of researchers, modular design and powerful infrastructure, exemplified by IBM’s Watson [39]. For this reason, much academic research focuses on settings in which the scope of the task is reduced. This has been achieved by restricting questions to a specific topic or domain [89], allowing systems access to pre-specified passages of text from which the answer can be inferred [58, 120], or centering both questions and answers on a particular knowledge base [13, 17].

In what follows, we show that the dictionary embedding models introduced in the previous sections may form a useful component of an open QA system. Given the absence of a knowledge base or web-scale information in our architecture, we narrow the scope of the task by focusing on general knowledge crossword questions. General knowledge (non-cryptic, or quick) crosswords appear in national newspapers in many countries. Crossword question answering is more tractable than general open QA for two reasons. First, models know the length of the correct answer (in letters), reducing the search space. Second, some crossword questions mirror definitions, in that they refer to fundamental properties of concepts (*a twelve-sided shape*) or request a category member (*a city in Egypt*).<sup>13</sup>

#### 4.4.1 Evaluation

General Knowledge crossword questions come in different styles and forms. We used the Eddie James crossword website to compile a bank of sentence-like general-knowledge questions.<sup>14</sup> Eddie James is one of the UK’s leading crossword compilers, working for several national newspapers. Our **long** question set consists of the first 150 questions (starting from puzzle #1) from his general-knowledge crosswords, excluding clues of fewer than four words and those whose answer was not a single word (e.g. *kingjames*).

To evaluate models on a different type of clue, we also compiled a set of **shorter** questions based on the Guardian Quick Crossword. Guardian questions still require general factual or linguistic knowledge, but are generally shorter and somewhat more cryptic than the longer Eddie James clues. We again formed a list of 150 questions, beginning on 1 January 2015 and excluding any questions with multiple-word answers.

---

<sup>13</sup>As our interest is in the language understanding, we do not address the question of fitting answers into a grid, which is the main concern of end-to-end automated crossword solvers [75].

<sup>14</sup><http://www.eddiejames.co.uk/>

For clear contrast, we excluded those few questions of length greater than four words. Of these 150 clues, a subset of 30 were **single-word** clues. All evaluation datasets are available online with the paper.

As with the reverse dictionary experiments, candidates are extracted from models by inputting definitions and returning words corresponding to the closest embeddings in the target space. In this case, however, we only consider candidate words *whose length matches the length specified in the clue*.

Test set	Word	Description
Long (150)	<i>Baudelaire</i>	"French poet and key figure in the development of Symbolism."
Short (120)	<i>satanist</i>	"devil devotee"
Single-Word (30)	<i>guilt</i>	"culpability"

**Table 4.5:** Examples of the different question types in the crossword question evaluation dataset.

Question Type		avg rank -accuracy@10/100 - rank variance							
		Long (150)			Short (120)			Single-Word (30)	
One Across		.39 /			<b>.68 /</b>			.70 /	
Crossword Maestro		.27 /			.43 /			.73 /	
W2V add		42	.31/.63	92	11	.50/.78	66	<b>2</b>	<b>.79/.90</b> 45
RNN cosine		15	.43/.69	108	22	.39/.67	117	72	.31/.52 187
RNN w2v cosine		4	.61/.82	60	<b>7</b>	.56/.79	60	12	.48/.72 116
RNN ranking		6	.58/.84	<b>48</b>	10	.51/.73	57	12	.48/.69 67
RNN w2v ranking		<b>3</b>	.62/.80	61	8	.57/.78	49	12	.48/.69 114
BOW cosine		4	.60/.82	54	<b>7</b>	.56/.78	51	12	.45/.72 137
BOW w2v cosine		4	.60/.83	56	<b>7</b>	.54/.80	48	3	.59/.79 111
BOW ranking		5	<b>.62/.87</b>	50	8	.58/. <b>83</b>	37	8	.55/.79 <b>39</b>
BOW w2v ranking		5	.60/.86	<b>48</b>	8	.56/.83	<b>35</b>	4	.55/.83 43

**Table 4.6:** Performance of different models on crossword questions of different length. The two commercial systems are evaluated via their web interface so only accuracy@10 can be reported in those cases.

## 4.4.2 Benchmarks and Comparisons

As with the reverse dictionary experiments, we compare RNN and BOW NLMs with a simple unsupervised baseline of elementwise addition of Word2Vec vectors in the em-

bedding space (we discard the ineffective *W2V mult* baseline), again restricting candidates to words of the pre-specified length. We also compare to two bespoke online crossword-solving engines. The first, One Across (<http://www.oneacross.com/>) is the candidate generation module of the award-winning *Proverb* crossword system [75]. *Proverb*, which was produced by academic researchers, has featured in national media such as *New Scientist*, and beaten expert humans in crossword solving tournaments. The second comparison is with Crossword Maestro (<http://www.crosswordmaestro.com/>), a commercial crossword solving system that handles both cryptic and non-cryptic crossword clues (we focus only on the non-cryptic setting), and has also been featured in national media.<sup>15</sup> We are unable to compare against a third well-known automatic crossword solver, *Dr Fill* [42], because code for *Dr Fill*'s candidate-generation module is not readily available. As with the RNN and baseline models, when evaluating existing systems we discard candidates whose length does not match the length specified in the clue.

Certain principles connect the design of the existing commercial systems and differentiate them from our approach. Unlike the NLMs, they each require query-time access to large databases containing common crossword clues, dictionary definitions, the frequency with which words typically appear as crossword solutions and other hand-engineered and task-specific components [75, 42].

### 4.4.3 Results

The performance of models on the various question types is presented in Table 4.6. When evaluating the two commercial systems, One Across and Crossword Maestro, we have access to web interfaces that return up to approximately 100 candidates for each query, so can only reliably record membership of the top ten (accuracy@10).

On the long questions, we observe a clear advantage for all dictionary embedding models over the commercial systems and the simple unsupervised baseline. Here, the best performing NLM (RNN with Word2Vec input embeddings and ranking loss) ranks the correct answer third on average, and in the top-ten candidates over 60% of the time.

As the questions get shorter, the advantage of the embedding models diminishes. Both the unsupervised baseline and One Across answer the short questions with comparable accuracy to the RNN and BOW models. One reason for this may be the dif-

---

<sup>15</sup> See e.g. <http://www.theguardian.com/crosswords/crossword-blog/2012/mar/08/crossword-blog-computers-crack-cryptic-clues>

Input Description	One Across	Crossword Maestro	BOW	RNN
"Swiss mountain peak famed for its north face (5)"	1: <i>noted</i> 2: <i>front</i> 3: <b>Eiger</b> 4: <i>crown</i> 5: <i>fount</i>	1: <i>after</i> 2: <i>favor</i> 3: <i>ahead</i> 4: <i>along</i> 5: <i>being</i>	1: <b>Eiger</b> 2: <i>Crags</i> 3: <i>Teton</i> 4: <i>Cerro</i> 5: <i>Jebel</i>	1: <b>Eiger</b> 2: <i>Aosta</i> 3: <i>Cuneo</i> 4: <i>Lecco</i> 5: <i>Tyrol</i>
"Old Testament successor to Moses (6)"	1: <b>Joshua</b> 2: <i>Exodus</i> 3: <i>Hebrew</i> 4: <i>person</i> 5: <i>across</i>	1: <i>devise</i> 2: <i>Daniel</i> 3: <i>Haggai</i> 4: <i>Isaiah</i> 5: <i>Joseph</i>	1: <i>Isaiah</i> 2: <i>Elijah</i> 3: <b>Joshua</b> 4: <i>Elisha</i> 5: <i>Yahweh</i>	1: <b>Joshua</b> 2: <i>Isaiah</i> 3: <i>Gideon</i> 4: <i>Elijah</i> 5: <i>Yahweh</i>
"The former currency of the Netherlands (7)"	1: <i>Holland</i> 2: <i>general</i> 3: <i>Lesotho</i>	1: <i>Holland</i> 2: <i>ancient</i> 3: <i>earlier</i> 4: <i>onetime</i> 5: <i>qondam</i>	1: <b>Guilder</b> 2: <i>Holland</i> 3: <i>Drenthe</i> 4: <i>Utrecht</i> 5: <i>Naarden</i>	1: <b>Guilder</b> 2: <i>Escudos</i> 3: <i>Pesetas</i> 4: <i>Someren</i> 5: <i>Florins</i>
"Arnold, 20th Century composer pioneer of atonality (10)"	1: <i>surrealism</i> 2: <i>laborparty</i> 3: <i>tonemusics</i> 4: <i>introduced</i> 5: <b>Schoenberg</b>	1: <i>disharmony</i> 2: <i>dissonance</i> 3: <i>bringabout</i> 4: <i>constitute</i> 5: <i>triggeroff</i>	1: <b>Schoenberg</b> 2: <i>Christleib</i> 3: <i>Stravinsky</i> 4: <i>Elderfield</i> 5: <i>Mendelsohn</i>	1: <i>Mendelsohn</i> 2: <i>Williamson</i> 3: <i>Huddleston</i> 4: <i>Mandelbaum</i> 5: <i>Zimmerman</i>

**Table 4.7:** Responses from different models to example crossword clues. In each case the model output is filtered to exclude any candidates that are not of the same length as the correct answer. BOW and RNN models are trained without Word2Vec input embeddings and cosine loss.

ference in form and style between the shorter clues and the full definitions or encyclopedia sentences in the dictionary training data. As the length of the clue decreases, finding the answer often reduces to generating synonyms (*culpability* - *guilt*), or category members (*tall animal* - *giraffe*). The commercial systems can retrieve good candidates for such clues among their databases of entities, relationships and common crossword answers. Unsupervised Word2Vec representations are also known to encode these sorts of relationships (even after elementwise addition for short sequences of words) [84]. This would also explain why the dictionary embedding models with pre-trained (Word2Vec) input embeddings outperform those with learned embeddings, particularly for the shortest questions.

#### 4.4.4 Qualitative Analysis

A better understanding of how the different models arrive at their answers can be gained from considering specific examples, as presented in Table 4.7. The first three examples show that, despite the apparently superficial nature of its training data (definitions and introductory sentences) embedding models can answer questions that require factual knowledge about people and places. Another notable characteristic of these

model is the consistent semantic appropriateness of the candidate set. In the first case, the top five candidates are all mountains, valleys or places in the Alps; in the second, they are all biblical names. In the third, the RNN model retrieves currencies, in this case performing better than the BOW model, which retrieves entities of various type associated with the Netherlands. Generally speaking (as can be observed by the web demo), the ‘smoothness’ or consistency in candidate generation of the dictionary embedding models is greater than that of the commercial systems. Despite its simplicity, the unsupervised W2V addition method is at times also surprisingly effective, as shown by the fact that it returns *Joshua* in its top candidates for the third query.

The final example in Table 4.7 illustrates the surprising power of the BOW model. In the training data there is a single definition for the correct answer *Schoenberg*: *United States composer and musical theorist (born in Austria) who developed atonal composition*. The only word common to both the query and the definition is ‘composer’ (there is no tokenization that allows the BOW model to directly connect *atonal* and *atonality*). Nevertheless, the model is able to infer the necessary connections between the concepts in the query and the definition to return Schoenberg as the top candidate.

Despite such cases, it remains an open question whether, with more diverse training data, the world knowledge required for full open QA (e.g. secondary facts about *Schoenberg*, such as his family) could be encoded and retained as weights in a (larger) dynamic network, or whether it will be necessary to combine the RNN with an external memory that is less frequently (or never) updated. This latter approach has begun to achieve impressive results on certain QA and entailment tasks [17, 45, 120].

## 4.5 Conclusion

Dictionaries exist in many of the world’s languages. We have shown how these lexical resources can constitute valuable data for training the latest neural language models to interpret and represent the meaning of phrases and sentences. While humans use the phrasal definitions in dictionaries to better understand the meaning of words, machines can use the words to better understand the phrases. We used two dictionary embedding architectures - a recurrent neural network architecture with a long-short-term memory, and a simpler linear bag-of-words model - to explicitly exploit this idea.

On the reverse dictionary task that mirrors its training setting, NLMs that embed

all known concepts in a continuous-valued vector space perform comparably to the best known commercial applications despite having access to many fewer definitions. Moreover, they generate smoother sets of candidates and require no linguistic pre-processing or task-specific engineering. We also showed how the description-to-word objective can be used to train models useful for other tasks. NLMs trained on the same data can answer general-knowledge crossword questions, and indeed outperform commercial systems on questions containing more than four words. While our QA experiments focused on crosswords, the results suggest that a similar embedding-based approach may ultimately lead to improved output from more general QA and dialog systems and information retrieval engines in general.

We make all code, training data, evaluation sets and both of our linguistic tools publicly available online for future research. In particular, we propose the reverse dictionary task as a comparatively general-purpose and objective way of evaluating how well models compose lexical meaning into phrase or sentence representations (whether or not they involve training on definitions directly).

In the next stage of this research, we will explore ways to enhance the NLMs described here, especially in the question-answering context. The models are currently not trained on any question-like language, and would conceivably improve on exposure to such linguistic forms. We would also like to understand better how BOW models can perform so well with no ‘awareness’ of word order, and whether there are specific linguistic contexts in which models like RNNs or others with the power to encode word order are indeed necessary. Finally, we intend to explore ways to endow the model with richer world knowledge. This may require the integration of an external memory module, similar to the promising approaches proposed in several recent papers [45, 120].

# Chapter 5

## Learning to Represent Sentences

### 5.0.1 Introduction

Distributed representations - dense real-valued vectors that encode the semantics of linguistic units - are ubiquitous in today's NLP research. For single-words or word-like entities, there are established ways to acquire such representations from naturally occurring (unlabelled) training data based on comparatively task-agnostic objectives (such as predicting adjacent words). These methods are well understood empirically [7] and theoretically [73]. The best word representation spaces reflect consistently-observed aspects of human conceptual organisation [54], and can be added as features to improve the performance of numerous language processing systems [29].

By contrast, there is comparatively little consensus on the best ways to learn distributed representations of phrases or sentences.<sup>1</sup> With the advent of deeper language processing techniques, it is relatively common for models to represent phrases or sentences as continuous-valued vectors. Examples include machine translation [110], image captioning [77] and dialogue systems [102]. While it has been observed informally that the internal sentence representations of such models can reflect semantic intuitions [25], it is not known which architectures or objectives yield the 'best' or most useful representations. Resolving this question could ultimately have a significant impact on language processing systems. Indeed, it is phrases and sentences, rather than individual words, that encode the human-like general world knowledge (or 'common sense') [93] that is a critical missing part of most current language understanding systems.

---

<sup>1</sup>See the contrasting conclusions in [86, 27, 6, 85] among others.

We address this issue with a systematic comparison of cutting-edge methods for learning distributed representations of sentences. We constrain our comparison to methods that do not require labelled data gathered for the purpose of training models, since such methods are more cost-effective and applicable across languages and domains. We also propose two new phrase or sentence representation learning objectives - *Sequential Denoising Autoencoders* (SDAEs) and *FastSent*, a sentence-level log-linear bag-of-words model. We compare all methods on two types of task - *supervised* and *unsupervised evaluations* - reflecting different ways in which representations are ultimately to be used. In the former setting, a classifier or regression model is applied to representations and trained with task-specific labelled data, while in the latter, representation spaces are directly queried using cosine distance.

We observe notable differences in approaches depending on the nature of the evaluation metric. In particular, deeper or more complex models (which require greater time and resources to train) generally perform best in the supervised setting, whereas shallow log-linear models work best on unsupervised benchmarks. Specifically, SkipThought Vectors [66] perform best on the majority of supervised evaluations, but SDAEs are the top performer on paraphrase identification. In contrast, on the (unsupervised) SICK sentence relatedness benchmark, FastSent, a simple, log-linear variant of the SkipThought objective, performs better than all other models. Interestingly, the method that exhibits strongest performance across both supervised and unsupervised benchmarks is a bag-of-words model trained to compose word embeddings using dictionary definitions [51]. Taken together, these findings constitute valuable guidelines for the application of phrasal or sentential representation-learning to language understanding systems.

## 5.0.2 Distributed Sentence Representations

To constrain the analysis, we compare neural language models that compute sentence representations from unlabelled, naturally-occurring data, as with the predominant methods for word representations.<sup>2</sup> Likewise, we do not focus on ‘bottom up’ models where phrase or sentence representations are built from fixed mathematical operations on word vectors (although we do consider a canonical case - see CBOW below); these were already compared by [85]. Most space is devoted to our novel approaches, and

---

<sup>2</sup>This excludes innovative supervised sentence-level architectures including [107, 64] and many others.



we refer the reader to the original papers for more details of existing models.

### 5.0.3 Existing Models Trained on Text

**SkipThought Vectors** For consecutive sentences  $S_{i-1}, S_i, S_{i+1}$  in some document, the **SkipThought** model [66] is trained to predict target sentences  $S_{i-1}$  and  $S_{i+1}$  given source sentence  $S_i$ . As with all *sequence-to-sequence* models, in training the source sentence is ‘encoded’ by a Recurrent Neural Network (RNN) (with Gated Recurrent uUnits [25]) and then ‘decoded’ into the two target sentences in turn. Importantly, because RNNs employ a single set of update weights at each time-step, both the encoder and decoder are sensitive to the order of words in the source sentence.

For each position in a target sentence  $S_t$ , the decoder computes a softmax distribution over the model’s vocabulary. The cost of a training example is the sum of the negative log-likelihood of each correct word in the target sentences  $S_{i-1}$  and  $S_{i+1}$ . This cost is backpropagated to train the encoder (and decoder), which, when trained, can map sequences of words to a single vector.

**ParagraphVector** [70] proposed two log-linear models of sentence representation. The **DBOW** model learns a vector  $s$  for every sentence  $S$  in the training corpus which, together with word embeddings  $v_w$ , define a softmax distribution optimised to predict words  $w \in S$  given  $S$ . The  $v_w$  are shared across all sentences in the corpus. In the **DM** model,  $k$ -grams of consecutive words  $\{w_i \dots w_{i+k} \in S\}$  are selected and  $s$  is combined with  $\{v_{w_i} \dots v_{w_{i+k}}\}$  to make a softmax prediction (parameterised by additional weights) of  $w_{i+k+1}$ .

We used the Gensim implementation,<sup>3</sup> treating each sentence in the training data as a ‘paragraph’ as suggested by the authors. During training, both DM and DBOW models store representations for every sentence (as well as word) in the training corpus. Even on large servers it was therefore only possible to train models with representation size 200, and DM models whose combination operation was averaging (rather than concatenation).

**Bottom-Up Methods** We train **CBOW** and **SkipGram** word embeddings [84] on the Books corpus, and compose by elementwise addition as proposed by [87].<sup>4</sup>

We also compare to **C-PHRASE** [98], an approach that exploits a (supervised) parser to infer distributed semantic representations based on a syntactic parse of sen-

<sup>3</sup><https://radimrehurek.com/gensim/>

<sup>4</sup>We also tried multiplication but this gave very poor results.

tences. C-PHRASE achieves state-of-the-art results for distributed representations on several evaluations used in this study.<sup>5</sup>

**Non-Distributed Baseline** We implement a **TFIDF BOW** model in which the representation of sentence  $S$  encodes the count in  $S$  of a set of feature-words weighted by their *tfidf* in  $C$ , the corpus. The feature-words are the 200,000 most common words in  $C$ .

#### 5.0.4 Models Trained on Structured Resources

The following models rely on (freely-available) data that has more structure than raw text.

**DictRep** [51] trained neural language models to map dictionary definitions to pre-trained word embeddings of the words defined by those definitions. They experimented with **BOW** and **RNN** (with LSTM) encoding architectures and variants in which the input word embeddings were either learned or pre-trained (**+embs.**) to match the target word embeddings. We implement their models using the available code and training data.<sup>6</sup>

**CaptionRep** Using the same overall architecture, we trained (**BOW** and **RNN**) models to map captions in the COCO dataset [24] to pre-trained vector representations of images. The image representations were encoded by a deep convolutional network [111] trained on the ILSVRC 2014 object recognition task [101]. Multi-modal distributed representations can be encoded by feeding test sentences forward through the trained model.

**NMT** We consider the sentence representations learned by neural MT models. These models have identical architecture to SkipThought, but are trained on sentence-aligned translated texts. We used a standard architecture [25] on all available **En-Fr** and **En-De** data from the 2015 Workshop on Statistical MT (WMT).<sup>7</sup>

---

<sup>5</sup>Since code for C-PHRASE is not publicly-available we use the available pre-trained model (<http://cllc.cimcc.unitn.it/composes/cphrase-vectors.html>). Note this model is trained on  $3\times$  more text than others in this study.

<sup>6</sup><https://www.cl.cam.ac.uk/~fh295/>. Definitions from the training data matching those in the WordNet STS 2014 evaluation (used in this study) were excluded.

<sup>7</sup>[www.statmt.org/wmt15/translation-task.html](http://www.statmt.org/wmt15/translation-task.html)

### 5.0.5 Novel Text-Based Models

We introduce two new approaches designed to address certain limitations with the existing models.

**Sequential (Denoising) Autoencoders** The SkipThought objective requires training text with a coherent inter-sentence narrative, making it problematic to port to domains such as social media or artificial language generated from symbolic knowledge. To avoid this restriction, we experiment with a representation-learning objective based on *denoising autoencoders* (DAEs). In a DAE, high-dimensional input data is corrupted according to some noise function, and the model is trained to recover the original data from the corrupted version. As a result of this process, DAEs learn to represent the data in terms of features that explain its important factors of variation [114]. Transforming data into DAE representations (as a ‘pre-training’ or initialisation step) gives more robust (supervised) classification performance in deep feedforward networks [115].

The original DAEs were feedforward nets applied to (image) data of fixed size. Here, we adapt the approach to variable-length sentences by means of a noise function  $N(S|p_o, p_x)$ , determined by free parameters  $p_o, p_x \in [0, 1]$ . First, for each word  $w$  in  $S$ ,  $N$  deletes  $w$  with (independent) probability  $p_o$ . Then, for each non-overlapping bigram  $w_i w_{i+1}$  in  $S$ ,  $N$  swaps  $w_i$  and  $w_{i+1}$  with probability  $p_x$ . We then train the same LSTM-based encoder-decoder architecture as NMT, but with the denoising objective to predict (as target) the original source sentence  $S$  given a corrupted version  $N(S|p_o, p_x)$  (as source). The trained model can then encode novel word sequences into distributed representations. We call this model the *Sequential Denoising Autoencoder* (**SDAE**). Note that, unlike SkipThought, SDAEs can be trained on sets of sentences in arbitrary order.

We label the case with no noise (i.e.  $p_o = p_x = 0$  and  $N \equiv id$ ) **SAE**. This setting matches the method applied to text classification tasks by [31]. The ‘word dropout’ effect when  $p_o \geq 0$  has also been used as a regulariser for deep nets in supervised language tasks [59], and for large  $p_x$  the objective is similar to word-level ‘debugging’ [109]. For the SDAE, we tuned  $p_o, p_x$  on the validation set (see Section 5.0.9).<sup>8</sup> We also tried a variant (**+embs**) in which words are represented by (fixed) pre-trained embeddings.

**FastSent** The performance of SkipThought vectors shows that rich sentence seman-

---

<sup>8</sup>We searched  $p_o, p_x \in \{0.1, 0.2, 0.3\}$  and observed best results with  $p_o = p_x = 0.1$ .

tics can be inferred from the content of adjacent sentences. The model could be said to exploit a type of *sentence-level Distributional Hypothesis* [48, 99]. Nevertheless, like many deep neural language models, SkipThought is very slow to train (see Table 5.1). FastSent is a simple additive (log-linear) sentence model designed to exploit the same signal, but at much lower computational expense. Given a BOW representation of some sentence in context, the model simply predicts adjacent sentences (also represented as BOW) .

More formally, FastSent learns a source  $u_w$  and target  $v_w$  embedding for each word in the model vocabulary. For a training example  $S_{i-1}, S_i, S_{i+1}$  of consecutive sentences,  $S_i$  is represented as the sum of its source embeddings  $\mathbf{s}_i = \sum_{w \in S_i} u_w$ . The cost of the example is then simply:

$$\sum_{w \in S_{i-1} \cup S_{i+1}} \phi(\mathbf{s}_i, v_w) \quad (5.1)$$

where  $\phi(v_1, v_2)$  is the softmax function.

We also experiment with a variant (**+AE**) in which the encoded (source) representation must predict its own words as target in addition to those of adjacent sentences. Thus in FastSent+AE, (5.1) becomes

$$\sum_{w \in S_{i-1} \cup S_i \cup S_{i+1}} \phi(\mathbf{s}_i, v_w). \quad (5.2)$$

At test time the trained model (very quickly) encodes unseen word sequences into distributed representations with  $\mathbf{s} = \sum_{w \in S} u_w$ .

### 5.0.6 Training and Model Selection

Unless stated above, all models were trained on the Toronto Books Corpus,<sup>9</sup> which has the inter-sentential coherence required for SkipThought and FastSent. The corpus consists of 70m ordered sentences from over 7,000 books.

Specifications of the models are shown in Table 5.1. The log-linear models (SkipGram, CBOW, ParagraphVec and FastSent) were trained for one epoch on one CPU core. The representation dimension  $d$  for these models was found after tuning  $d \in \{100, 200, 300, 400, 500\}$  on the validation set.<sup>10</sup> All other models were trained on one

<sup>9</sup><http://www.cs.toronto.edu/~mbweb/>

<sup>10</sup>For ParagraphVec only  $d \in \{100, 200\}$  was possible due to the high memory footprint.

	OS	R	WO	SD	WD	TR	TE
S(D)AE			✓	2400	100	72*	640
ParagraphVec				100	100	4	1130
CBOW				500	500	2	145
SkipThought	✓		✓	4800	620	336*	890
FastSent	✓			100	100	2	140
DictRep		✓	✓	500	256	24*	470
CaptionRep		✓	✓	500	256	24*	470
NMT		✓	✓	2400	512	72*	720

**Table 5.1: Properties of models compared in this study** **OS:** requires training corpus of sentences in order. **R:** requires structured resource for training. **WO:** encoder sensitive to word order. **SD:** dimension of sentence representation. **WD:** dimension of word representation. **TR:** approximate training time (hours) on the dataset in this paper. \* indicates trained on GPU. **TE:** approximate time (s) taken to encode 0.5m sentences.

Dataset	Sentence 1	Sentence 2
News	<i>Mexico wishes to guarantee citizens' safety.</i>	<i>Mexico wishes to avoid more violence.</i>
Forum	<i>The problem is simpler than that.</i>	<i>The problem is simple.</i>
STS WordNet	<i>A social set or clique of friends.</i>	<i>An unofficial association of people or groups.</i>
2014 Twitter	<i>Taking Aim #Stopgunviolence #Congress #NRA</i>	<i>Obama, Gun Policy and the N.R.A.</i>
Images	<i>A woman riding a brown horse.</i>	<i>A young girl riding a brown horse.</i>
Headlines	<i>Iranians Vote in Presidential Election.</i>	<i>Keita Wins Mali Presidential Election.</i>
SICK (test+train)	<i>A lone biker is jumping in the air.</i>	<i>A man is jumping into a full pool.</i>

**Table 5.2:** Example sentence pairs and ‘similarity’ ratings from the unsupervised evaluations used in this study.

GPU. The S(D)AE models were trained for one epoch ( $\approx 8$  days). The SkipThought model was trained for two weeks, covering just under one epoch.<sup>11</sup> For CaptionRep and DictRep, performance was monitored on held-out training data and training was stopped after 24 hours after a plateau in cost. The NMT models were trained for 72 hours.

### 5.0.7 Evaluating Sentence Representations

In previous work, distributed representations of language were evaluated either by measuring the effect of adding representations as features in some classification task - *supervised evaluation* [29, 81, 66] - or by comparing with human relatedness judgments - *unsupervised evaluation* [51, 7, 74]. The former setting reflects a scenario in

<sup>11</sup>Downloaded from <https://github.com/ryankiros/skip-thoughts>

which representations are used to inject general knowledge (sometimes considered as *pre-training*) into a supervised model. The latter pertains to applications in which the sentence representation space is used for direct comparisons, lookup or retrieval. Here, we apply and compare both evaluation paradigms.

Data	Model	MSRP (Acc / F1)	MR	CR	SUBJ	MPQA	TREC
Unordered Sentences (Toronto Books: 70m sents, 0.9B words)	SAE	74.3 / 81.7	62.6	68.0	86.1	76.8	80.2
	SAE+embs.	70.6 / 77.9	73.2	75.3	89.8	86.2	80.4
	SDAE	<b><u>76.4 / 83.4</u></b>	67.6	74.0	89.3	81.3	77.6
	SDAE+embs.	73.7 / 80.7	<b>74.6</b>	<b>78.0</b>	<b>90.8</b>	<b>86.9</b>	78.4
	ParagraphVec DBOW	72.9 / 81.1	60.2	66.9	76.3	70.7	59.4
	ParagraphVec DM	73.6 / 81.9	61.5	68.6	76.4	78.1	55.8
	Skipgram	69.3 / 77.2	73.6	77.3	89.2	85.0	82.2
	CBOW	67.6 / 76.1	73.6	77.3	89.1	85.0	82.2
	Unigram TFIDF	<b>73.6 / 81.7</b>	73.7	79.2	90.3	82.4	<b>85.0</b>
Ordered Sentences (Toronto Books)	SkipThought	<b>73.0 / 82.0</b>	<b>76.5</b>	<b>80.1</b>	<b>93.6</b>	<b>87.1</b>	<b>92.2</b>
	FastSent	72.2 / 80.3	70.8	78.4	88.7	80.6	76.8
	FastSent+AE	71.2 / 79.1	71.8	76.7	88.8	81.5	80.4
Other structured data resource	NMT En to Fr	69.1 / 77.1	64.7	70.1	84.9	81.5	<b>82.8</b>
	NMT En to De	65.2 / 73.3	61.0	67.6	78.2	72.9	81.6
	CaptionRep BOW	73.6 / 81.9	61.9	69.3	77.4	70.8	72.2
	CaptionRep RNN	72.6 / 81.1	55.0	64.9	64.9	71.0	62.4
	DictRep BOW	<b>73.7 / 81.6</b>	71.3	75.6	86.6	82.5	73.8
	DictRep BOW+embs.	68.4 / 76.8	<b>76.7</b>	<b>78.7</b>	<b>90.7</b>	<b>87.2</b>	81.0
	DictRep RNN	73.2 / 81.6	67.8	72.7	81.4	82.5	75.8
	DictRep RNN+embs.	66.8 / 76.0	72.5	73.5	85.6	85.7	72.0
2.8B words	CPHRASE	72.2 / 79.6	75.7	78.8	91.1	86.2	78.8

**Table 5.3:** Performance of sentence representation models on **supervised** evaluations (Section 5.0.8). Bold numbers indicate best performance in class. Underlined indicates best overall.

## 5.0.8 Supervised Evaluations

Representations are applied to 6 sentence classification tasks: paraphrase identification (MSRP) [34], movie review sentiment (MR) [96], product reviews (CR) [56], subjectivity classification (SUBJ) [95], opinion polarity (MPQA) [121] and question type classification (TREC) [117]. We follow the procedure (and code) of [66]: a logistic regression classifier is trained on top of sentence representations, with 10-fold cross-validation used when a train-test split is not pre-defined.

Model	STS 2014							SICK Test + Train
	News	Forum	WordNet	Twitter	Images	Headlines	All	
SAE	17/.16	.12/.12	.30/.23	.28/.22	.49/.46	.13/.11	.12/.13	.32/.31
SAE+embs.	.52/.54	.22/.23	.60/.55	.60/.60	.64/.64	.41/.41	.42/.43	.47/.49
SDAE	.07/.04	.11/.13	.33/.24	.44/.42	.44/.38	.36/.36	.17/.15	.46/.46
SDAE+embs.	.51/.54	.29/.29	.56/.50	.57/.58	.59/.59	.43/.44	.37/.38	.46/.46
ParagraphVec DBOW	.31/.34	.32/.32	.53/.5	.43/.46	.46/.44	.39/.41	.42/.43	.42/.46
ParagraphVec DM	.42/.46	.33/.34	.51/.48	.54/.57	.32/.30	.46/.47	.44/.44	.44/.46
Skipgram	.56/.59	.42/.42	<b>.73/.70</b>	<b>.71/.74</b>	.65/.67	<b>.55/.58</b>	.62/.63	<b>.60/.69</b>
CBOW	<b>.57/.61</b>	<b>.43/.44</b>	.72/.69	<b>.71/.75</b>	.71/.73	<b>.55/.59</b>	<b>.64/.65</b>	<b>.60/.69</b>
Unigram TFIDF	.48/.48	.40/.38	.60/.59	.63/.65	<b>.72/.74</b>	.49/.49	.58/.57	.52/.58
SkipThought	.44/.45	.14/.15	.39/.34	.42/.43	.55/.60	.43/.44	.27/.29	.57/.60
FastSent	<b>.58/.59</b>	<b>.41/.36</b>	<b>.74/.70</b>	.63/.66	<b>.74/.78</b>	.57/.59	<b>.63/.64</b>	<b>.61/.72</b>
FastSent+AE	.56/ <b>.59</b>	<b>.41/.40</b>	.69/.64	<b>.70/.74</b>	.63/.65	<b>.58/.60</b>	.62/.62	.60/.65
NMT En to Fr	.35/.32	.18/.18	.47/.43	.55/.53	.44/.45	.43/.43	.43/.42	.47/.49
NMT En to De	.47/.43	.26/.25	.34/.31	.49/.45	.44/.43	.38/.37	.40/.38	.46/.46
CaptionRep BOW	.26/.26	.29/.22	.50/.35	.37/.31	<b>.78/.81</b>	.39/.36	.46/.42	.56/.65
CaptionRep RNN	.05/.05	.13/.09	.40/.33	.36/.30	<b>.76/.82</b>	.30/.28	.39/.36	.53/.62
DictRep BOW	.62/.67	.42/.40	.81/.81	.62/.66	.66/.68	.53/.58	.62/.65	.57/.66
DictRep BOW+embs.	<b>.65/.72</b>	<b>.49/.47</b>	<b>.85/.86</b>	<b>.67/.72</b>	.71/.74	<b>.57/.61</b>	<b>.67/.70</b>	<b>.61/.70</b>
DictRep RNN	.40/.46	.26/.23	.78/.78	.42/.42	.56/.56	.38/.40	.49/.50	.49/.56
DictRep RNN+embs.	.51/.60	.29/.27	.80/.81	.44/.47	.65/.70	.42/.46	.54/.57	.49/.59
CPHRASE	<b>.69/.71</b>	.43/.41	.76/.73	.60/.65	.75/.79	<b>.60/.65</b>	.65/.67	<b>.60/.72</b>

**Table 5.4:** Performance of sentence representation models (Spearman/Pearson correlations) on **unsupervised** (relatedness) evaluations (Section 5.0.9). Models are grouped according to training data as indicated in Table 5.3.

## 5.0.9 Unsupervised Evaluations

We also measure how well representation spaces reflect human intuitions of the semantic sentence relatedness, by computing the cosine distance between vectors for the two sentences in each test pair, and correlating these distances with gold-standard human judgements. The SICK dataset [78] consists of 10,000 pairs of sentences and relatedness judgements. The STS 2014 dataset [2] consists of 3,750 pairs and ratings from six linguistic domains. Example ratings are shown in Table 5.2. All available pairs are used for testing apart from the 500 SICK ‘trial’ pairs, which are held-out for tuning hyperparameters (representation size of log-linear models, and noise parameters in SDAE). The optimal settings on this task are then applied to both supervised and unsupervised evaluations.

### 5.0.10 Results

Performance of the models on the supervised evaluations (grouped according to the data required by their objective) is shown in Table 5.3. Overall, SkipThought vectors perform best on three of the six evaluations, the BOW DictRep model with pre-trained word embeddings performs best on two, and the SDAE on one. SDAEs perform notably well on the paraphrasing task, going beyond SkipThought by three percentage points and approaching state-of-the-art performance of models designed specifically for the task [61]. SDAE is also consistently better than SAE, which aligns with other findings that adding noise to AEs produces richer representations [114].

Results on the unsupervised evaluations are shown in Table 5.4. The same DictRep model performs best on four of the six STS categories (and overall) and is joint-top performer on SICK. Of the models trained on raw text, simply adding CBOW word vectors works best on STS. The best performing raw text model on SICK is FastSent, which achieves almost identical performance to C-PHRASE’s state-of-the-art performance for a distributed model [98]. Further, it uses less than a third of the training text and does not require access to (supervised) syntactic representations for training. Together, the results of FastSent on the unsupervised evaluations and SkipThought on the supervised benchmarks provide strong support for the sentence-level distributional hypothesis: the context in which a sentence occurs provides valuable information about its semantics.

Across both unsupervised and supervised evaluations, the BOW DictRep with pre-trained word embeddings exhibits by some margin the most consistent performance. This robust performance suggests that DictRep representations may be particularly valuable when the ultimate application is non-specific or unknown, and confirms that dictionary definitions (where available) can be a powerful resource for representation learning.

### 5.0.11 Discussion

Many additional conclusions can be drawn from the results in Tables 5.3 and 5.4.

**Different objectives yield different representations** It may seem obvious, but the results confirm that different learning methods are preferable for different intended applications (and this variation appears greater than for word representations). For instance, it is perhaps unsurprising that SkipThought performs best on TREC because



the labels in this dataset are determined by the language immediately following the represented question (i.e. the answer) [117]. Paraphrase detection, on the other hand, may be better served by a model that focused entirely on the content *within* a sentence, such as SDAEs. Similar variation can be observed in the unsupervised evaluations. For instance, the (multimodal) representations produced by the CaptionRep model do not perform particularly well apart from on the Image category of STS where they beat all other models, demonstrating a clear effect of the well-studied modality differences in representation learning [20].

The nearest neighbours in Table 5.5 give a more concrete sense of the representation spaces. One notable difference is between (AE-style) models whose semantics come from within-sentence relationships (CBOW, SDAE, DictRep, ParagraphVec) and SkipThought/FastSent, which exploit the context around sentences. In the former case, nearby sentences generally have a high proportion of words in common, whereas for the latter it is the general concepts and/or function of the sentence that is similar, and word overlap is often minimal. Indeed, this may be a more important trait of FastSent than the marginal improvement on the SICK task. Readers can compare the CBOW and FastSent spaces at <http://45.55.60.98/>.

**Differences between supervised and unsupervised performance** Many of the best performing models on the supervised evaluations do not perform well in the unsupervised setting. In the SkipThought, S(D)AE and NMT models, the cost is computed based on a non-linear decoding of the internal sentence representations, so, as also observed by [3], the informative geometry of the representation space may not be reflected in a simple cosine distance. The log-linear models generally perform better in this unsupervised setting.

**Differences in resource requirements** As shown in Table 5.1, different models require different resources to train and use. This can limit their possible applications. For instance, while it was easy to make an online demo for fast querying of near neighbours in the CBOW and FastSent spaces, it was not practical for other models owing to memory footprint, encoding time and representation dimension.

**The role of word order is unclear** The average scores of models that are sensitive to word order (76.3) and of those that are not (76.6) are approximately the same across

Query	<i>If he had a weapon, he could maybe take out their last imp, and then beat up Errol and Vanessa.</i>	<i>An annoying buzz started to ring in my ears, becoming louder and louder as my vision began to swim.</i>
CBOW	<i>Then Rob and I would duke it out, and every once in a while, he would actually beat me.</i>	<i>Louder.</i>
Skip Thought	<i>If he could ram them from behind, send them sailing over the far side of the levee, he had a chance of stopping them.</i>	<i>A weighty pressure landed on my lungs and my vision at the edges, threatening my consciousness altogether.</i>
FastSent	<i>Isak's close enough to pick off any one of them, maybe all of them, if he had his rifle and a mind to.</i>	<i>The noise grew louder, the quaking increased as the sidewalk beneath my feet began to tremble even more.</i>
SDAE	<i>He'd even killed some of the most dangerous criminals in the galaxy, but none of those men had gotten to him like Vitktis.</i>	<i>I smile because I'm familiar with the knock, pausing to take a deep breath before dashing down.</i>
DictRep (FF+embs.)	<i>Kevin put a gun to the man's head, but even though he cried, he couldn't tell Kevin anything more.</i>	<i>Then gradually I began to hear a ringing in my ears.</i>
Paragraph Vector (DM)	<i>I take a deep breath and open the doors.</i>	<i>They listened as the motorcycle-like roar of an engine got louder and louder then stopped.</i>

**Table 5.5:** Sample nearest neighbour queries selected from a randomly sampled 0.5m sentences of the Toronto Books Corpus.

Supervised (combined $\alpha = 0.90$ )						Unsupervised (combined $\alpha = 0.93$ )						
MSRP	MR	CR	SUBJ	MPAQ	TREC	News	Forum	WordNet	Twitter	Images	Headlines	All STS
0.94 (6)	0.85 (1)	0.86 (4)	0.85 (1)	0.86 (3)	0.89 (5)	0.92 (4)	0.92 (3)	0.92 (4)	0.93 (6)	0.95 (8)	0.92 (2)	0.91 (1)

**Table 5.6:** Internal consistency (Chronbach's  $\alpha$ ) among evaluations when individual benchmarks are left out of the (supervised or unsupervised) cohorts. Consistency rank within cohort is in parentheses (1 = most consistent with other evaluations).

supervised evaluations. Across the unsupervised evaluations, however, BOW models score 0.55 on average compared with 0.42 for RNN-based (order sensitive) models. This seems at odds with the widely held view that word order plays an important role in determining the meaning of English sentences. One possibility is that order-critical sentences that cannot be disambiguated by a robust conceptual semantics (that could be encoded in distributed lexical representations) are in fact relatively rare. However, it is also plausible that current available evaluations do not adequately reflect order-dependent aspects of meaning (see below). This latter conjecture is supported by the comparatively strong performance of TFIDF BOW vectors, in which the effective lexical semantics are limited to simple relative frequencies.

**The evaluations have limitations** The internal consistency (Chronbach's  $\alpha$ ) of all evaluations considered together is 0.81 (just above 'acceptable').<sup>12</sup> Table 5.6 shows that consistency is far higher ('excellent') when considering the supervised or unsupervised tasks as independent cohorts. This indicates that, with respect to common characteristics of sentence representations, the supervised and unsupervised benchmarks do indeed prioritise different properties. It is also interesting that, by this metric, the properties measured by MSRP and image-caption relatedness are the furthest

<sup>12</sup>[wikipedia.org/wiki/Cronbach's\\_alpha](https://wikipedia.org/wiki/Cronbach's_alpha)

removed from other evaluations in their respective cohorts.

While these consistency scores are a promising sign, they could also be symptomatic of a set of evaluations that are all limited in the same way. The inter-rater agreement is only reported for one of the 8 evaluations considered (MPQA, 0.72 [121]), and for MR, SUBJ and TREC, each item is only rated by one or two annotators to maximise coverage. Table 5.2 illustrates why this may be an issue for the unsupervised evaluations; the notion of sentential 'relatedness' seems very subjective. It should be emphasised, however, that the tasks considered in this study are all frequently used for evaluation, and, to our knowledge, there are no existing benchmarks that overcome these limitations.

### 5.0.12 Conclusion

Advances in deep learning algorithms, software and hardware mean that many architectures and objectives for learning distributed sentence representations from unlabelled data are now available to NLP researchers. We have presented the first (to our knowledge) systematic comparison of these methods. We showed notable variation in the performance of approaches across a range of evaluations. Among other conclusions, we found that the optimal approach depends critically on whether representations will be applied in supervised or unsupervised settings - in the latter case, fast, shallow BOW models can still achieve the best performance. Further, we proposed two new objectives, FastSent and Sequential Denoising Autoencoders, which perform particularly well on specific tasks (MSRP and SICK sentence relatedness respectively).<sup>13</sup> If the application is unknown, however, the best all round choice may be DictRep: learning a mapping of pre-trained word embeddings from the word-phrase signal in dictionary definitions. While we have focused on models using naturally-occurring training data, in future work we will also consider supervised architectures (including convolutional, recursive and character-level models), potentially training them on multiple supervised tasks as an alternative way to induce the 'general knowledge' needed to give language technology the elusive human touch.

---

<sup>13</sup>We make all code for training and evaluating these new models publicly available, together with pre-trained models and an online demo of the FastSent sentence space.



## **Chapter 6**

# **Representation Learning**

Some representation learning text



# **Chapter 7**

## **Conclusion**

Some conclusive learning text





# Bibliography

- [1] Eneko Agirre, Enrique Alfonseca, Keith Hall, Jana Kravalova, Marius Pasca, and Aitor Soroa. A study on similarity and relatedness using distributional and wordnet-based approaches. In *Proceedings of NAACL-HLT 2009*, 2009.
- [2] Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Rada Mihalcea, German Rigau, and Janyce Wiebe. Semeval-2014 task 10: Multilingual semantic textual similarity. In *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, pages 81–91, 2014.
- [3] Amjad Almahairi, Kyle Kastner, Kyunghyun Cho, and Aaron Courville. Learning distributed representations from reviews for collaborative filtering. In *Proceedings of the 9th ACM Conference on Recommender Systems*, pages 147–154. ACM, 2015.
- [4] Mark Andrews, Gabriella Vigliocco, and David Vinson. Integrating experiential and distributional data to learn semantic representations. *Psychological Review*, 116(3):463, 2009.
- [5] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *Proceeding of ICLR*, 2015.
- [6] Marco Baroni, Raffaella Bernardi, and Roberto Zamparelli. Frege in space: A program of compositional distributional semantics. *Linguistic Issues in Language Technology*, 9, 2014.
- [7] Marco Baroni, Georgiana Dinu, and Germán Kruszewski. 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, 2014.
- [8] Lawrence W Barsalou. Grounded cognition: past, present, and future. *Topics in Cognitive Science*, 2(4):716–724, 2010.
  - [9] Lawrence W Barsalou and Katja Wiemer-Hastings. Situating abstract concepts. *Grounding Cognition: The Role of Perception and Action in Memory, Language, and Thought*, pages 129–163, 2005.
  - [10] Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. A neural probabilistic language model. *J. Mach. Learn. Res.*, 3:1137–1155, March 2003.
  - [11] Yoshua Bengio and Jean-Sébastien S  n  cal. Quick training of probabilistic neural nets by importance sampling. In *Proceedings of AISTATS 2003*, 2003.
  - [12] Yoshua Bengio, Patrice Simard, and Paolo Frasconi. Learning long-term dependencies with gradient descent is difficult. *Neural Networks, IEEE Transactions on*, 5(2):157–166, 1994.
  - [13] Jonathan Berant and Percy Liang. Semantic parsing via paraphrasing. In *Proceedings of the Association for Computational Linguistics*, 2014.
  - [14] James Bergstra, Olivier Breuleux, Fr  d  ric Bastien, Pascal Lamblin, Razvan Pascanu, Guillaume Desjardins, Joseph Turian, David Warde-Farley, and Yoshua Bengio. Theano: a CPU and GPU math expression compiler. In *Proceedings of the Python for Scientific Computing Conference (SciPy)*, June 2010. Oral Presentation.
  - [15] Slaven Bilac, Timothy Baldwin, and Hozumi Tanaka. Improving dictionary accessibility by maximizing use of available knowledge. *Traitement Automatique des Langues*, 44(2):199–224, 2003.
  - [16] Slaven Bilac, Wataru Watanabe, Taiichi Hashimoto, Takenobu Tokunaga, and Hozumi Tanaka. Dictionary search based on the target word description. In *Proceedings of NLP 2014*, 2004.
  - [17] Antoine Bordes, Sumit Chopra, and Jason Weston. Question answering with subgraph embeddings. *Proceedings of EMNLP*, 2014.

- [18] Antoine Bordes, Nicolas Usunier, Sumit Chopra, and Jason Weston. Large-scale simple question answering with memory networks. *arXiv preprint arXiv:1506.02075*, 2015.
- [19] Elia Bruni, Gemma Boleda, Marco Baroni, and Nam-Khanh Tran. Distributional semantics in technicolor. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1*, pages 136–145. Association for Computational Linguistics, 2012.
- [20] Elia Bruni, Nam Khanh Tran, and Marco Baroni. Multimodal distributional semantics. *Journal of Artificial Intelligence Research*, 49:1–47, 2014.
- [21] Joan L Bybee and Paul J Hopper. *Frequency and the Emergence of Linguistic Structure*, volume 45. John Benjamins Publishing, 2001.
- [22] Sarath Chandar, Stanislas Lauly, Hugo Larochelle, Mitesh Khapra, Balaraman Ravindran, Vikas C. Raykar, and Amrita Saha. An autoencoder approach to learning bilingual word representations. In *Advances in Neural Information Processing Systems*, pages 1853–1861, 2014.
- [23] Nick Chater and Christopher D Manning. Probabilistic models of language processing and acquisition. *Trends in Cognitive Sciences*, 10(7):335–344, 2006.
- [24] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollar, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015.
- [25] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *Proceedings of EMNLP*, 2014.
- [26] Kyunghyun Cho, Bart van Merrienboer, Caglar Gulcehre, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of EMNLP*, 2014.
- [27] Stephen Clark and Stephen Pulman. Combining symbolic and distributional models of meaning. In *AAAI Spring Symposium: Quantum Interaction*, pages 52–55, 2007.

- [28] Ronan Collobert and Jason Weston. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, pages 160–167. ACM, 2008.
- [29] Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. Natural language processing (almost) from scratch. *The Journal of Machine Learning Research*, 12:2493–2537, 2011.
- [30] Sebastian J Crutch and Elizabeth K Warrington. Abstract and concrete concepts have structurally different representational frameworks. *Brain*, 128(3):615–627, 2005.
- [31] Andrew M Dai and Quoc V Le. Semi-supervised sequence learning. In *Advances in Neural Information Processing Systems*, pages 3061–3069, 2015.
- [32] Barry J Devereux, Lorraine K Tyler, Jeroen Geertzen, and Billi Randall. The centre for speech, language and the brain (cslb) concept property norms. *Behavior Research Methods*, pages 1–9, 2013.
- [33] Jacob Devlin, Rabih Zbib, Zhongqiang Huang, Thomas Lamar, Richard Schwartz, and John Makhoul. Fast and robust neural network joint models for statistical machine translation. In *Proceedings of the Association for Computational Linguistics*, 2014.
- [34] Bill Dolan, Chris Quirk, and Chris Brockett. Unsupervised construction of large paraphrase corpora: Exploiting massively parallel news sources. In *Proceedings of the 20th international conference on Computational Linguistics*, page 350. Association for Computational Linguistics, 2004.
- [35] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *The Journal of Machine Learning Research*, 12:2121–2159, 2011.
- [36] Manaal Faruqui, Jesse Dodge, Sujay K. Jauhar, Chris Dyer, Eduard Hovy, and Noah A. Smith. Retrofitting word vectors to semantic lexicons. *Proceedings of the North American Chapter of the Association for Computational Linguistics*, 2014.

- [37] Manaal Faruqui and Chris Dyer. Improving vector space word representations using multilingual correlation. In *Proceedings of EACL*, volume 2014, 2014.
- [38] Yansong Feng and Mirella Lapata. Visual information in semantic representation. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 91–99. Association for Computational Linguistics, 2010.
- [39] David Ferrucci, Eric Brown, Jennifer Chu-Carroll, James Fan, David Gondek, Aditya A. Kalyanpur, Adam Lally, J. William Murdock, Eric Nyberg, John Prager, Nico Schlaefter, and Chris Welty. Building Watson: An overview of the DeepQA project. In *AI magazine*, volume 31(3), pages 59–79, 2010.
- [40] R. Firth, J. *A synopsis of linguistic theory 1930-1955*, pages 1–32. Oxford: Philological Society, 1957.
- [41] Jerry A Fodor. *The modularity of mind: An essay on faculty psychology*. MIT press, 1983.
- [42] Matthew L. Ginsberg. Dr. FILL: Crosswords and an implemented solver for singly weighted CSPs. In *Journal of Artificial Intelligence Research*, pages 851–886, 2011.
- [43] Stephan Gouws, Yoshua Bengio, and Greg Corrado. BilBOWA: Fast bilingual distributed representations without word alignments. In *Proceedings of NIPS Deep Learning Workshop*, 2014.
- [44] Scott T Grafton. Embodied cognition and the simulation of action to understand others. *Annals of the New York Academy of Sciences*, 1156(1):97–117, 2009.
- [45] Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines. *arXiv preprint arXiv:1410.5401*, 2014.
- [46] Aria Haghighi, Percy Liang, Taylor Berg-Kirkpatrick, and Dan Klein. Learning bilingual lexicons from monolingual corpora. In *ACL*, volume 2008, pages 771–779, 2008.
- [47] David R Hardoon, Sandor Szedmak, and John Shawe-Taylor. Canonical correlation analysis: An overview with application to learning methods. *Neural Computation*, 16(12):2639–2664, 2004.

- [48] Zellig S Harris. Distributional structure. *Word*, 1954.
- [49] Karl Moritz Hermann and Phil Blunsom. Multilingual distributed representations without word alignment. In *Proceedings of ICLR*, 2013.
- [50] Karl Moritz Hermann and Phil Blunsom. Multilingual Distributed Representations without Word Alignment. In *Proceedings of ICLR*, 2014.
- [51] Felix Hill, Kyunghyun Cho, Anna Korhonen, and Yoshua Bengio. Learning to understand phrases by embedding the dictionary. *Transactions of the Association for Computational Linguistics*, 2015.
- [52] Felix Hill and Anna Korhonen. Learning abstract concepts from multi-modal data: Since you probably can’t see what i mean. In *Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014)*, October 2014.
- [53] Felix Hill, Anna Korhonen, and Christian Bentz. A quantitative empirical analysis of the abstract/concrete distinction. *Cognitive Science*, 2013.
- [54] Felix Hill, Roi Reichart, and Anna Korhonen. Simlex-999: Evaluating semantic models with (genuine) similarity estimation. *Computational Linguistics*, 2015.
- [55] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [56] Minqing Hu and Bing Liu. Mining and summarizing customer reviews. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 168–177. ACM, 2004.
- [57] Eric H Huang, Richard Socher, Christopher D Manning, and Andrew Y Ng. 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*, pages 873–882. Association for Computational Linguistics, 2012.
- [58] Mohit Iyyer, Jordan Boyd-Graber, Leonardo Claudino, Richard Socher, and Hal Daumé III. A neural network for factoid question answering over paragraphs. In *Proceedings of EMNLP*, 2014.

- [59] Mohit Iyyer, Varun Manjunatha, Jordan Boyd-Graber, and Hal Daumé III. Deep unordered composition rivals syntactic methods for text classification. In *Proceedings of the Association for Computational Linguistics*, 2015.
- [60] Sébastien Jean, Kyunghyun Cho, Roland Memisevic, and Yoshua Bengio. On using very large target vocabulary for neural machine translation. *Proceedings of ALC*, 2015.
- [61] Yangfeng Ji and Jacob Eisenstein. Discriminative improvements to distributional sentence similarity. In *EMNLP*, pages 891–896, 2013.
- [62] Brendan T Johns and Michael N Jones. Perceptual inference through global lexical similarity. *Topics in Cognitive Science*, 4(1):103–120, 2012.
- [63] Nal Kalchbrenner and Phil Blunsom. Recurrent continuous translation models. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, Seattle, October 2013. Association for Computational Linguistics.
- [64] Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. A convolutional neural network for modelling sentences. In *Proceedings of EMNLP*, 2014.
- [65] Ryan Kiros, Ruslan Salakhutdinov, and Richard S. Zemel. Unifying visual-semantic embeddings with multimodal neural language models. *Transactions of the Association for Computational Linguistics*, 2015. to appear.
- [66] Ryan Kiros, Yukun Zhu, Ruslan R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Skip-thought vectors. In *Advances in Neural Information Processing Systems*, pages 3276–3284, 2015.
- [67] A. Klementiev, I. Titov, and B. Bhattarai. Inducing crosslingual distributed representations of words. In *Proceedings of COLING*, 2012.
- [68] Tomáš Kočiský, Karl Moritz Hermann, and Phil Blunsom. Learning Bilingual Word Representations by Marginalizing Alignments. In *Proceedings of ACL*, jun 2014.
- [69] Thomas K Landauer and Susan T Dumais. A solution to plato’s problem: The latent semantic analysis theory of acquisition, induction, and representation of knowledge. *Psychological Review*, 104(2):211, 1997.

- [70] Quoc V Le and Tomas Mikolov. Distributed representations of sentences and documents. In *Proceedings of ICML*, 2014.
- [71] Geoffrey Leech, Roger Garside, and Michael Bryant. Claws4: the tagging of the British National Corpus. In *Proceedings of the 15th conference on Computational linguistics-Volume 1*, pages 622–628. Association for Computational Linguistics, 1994.
- [72] Omer Levy and Yoav Goldberg. Dependency-based word embeddings. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, volume 2, 2014.
- [73] Omer Levy and Yoav Goldberg. Neural word embedding as implicit matrix factorization. In *Advances in Neural Information Processing Systems*, pages 2177–2185, 2014.
- [74] Omer Levy, Yoav Goldberg, and Ido Dagan. Improving distributional similarity with lessons learned from word embeddings. *Transactions of the Association for Computational Linguistics*, 3:211–225, 2015.
- [75] Michael L. Littman, Greg A. Keim, and Noam Shazeer. A probabilistic approach to solving crossword puzzles. *Artificial Intelligence*, 134(1):23–55, 2002.
- [76] Thang Luong, Ilya Sutskever, Quoc V Le, Oriol Vinyals, and Wojciech Zaremba. Addressing the rare word problem in neural machine translation. *Proceedings of the Association for Computational Linguistics*, 2015.
- [77] Junhua Mao, Wei Xu, Yi Yang, Jiang Wang, and Alan Yulle. Deep captioning with multimodal recurrent neural networks (m-rnn). In *Proceedings of ICLR*, 2015.
- [78] Marco Marelli, Stefano Menini, Marco Baroni, Luisa Bentivogli, Raffaella Bernardi, and Roberto Zamparelli. A sick cure for the evaluation of compositional distributional semantic models. In *Proceedings of LREC*, pages 216–223. Citeseer, 2014.
- [79] Ken McRae, George S Cree, Mark S Seidenberg, and Chris McNorgan. Semantic feature production norms for a large set of living and nonliving things. *Behavior Research Methods*, 37(4):547–559, 2005.



- [80] Grégoire Mesnil, Yann Dauphin, Xavier Glorot, Salah Rifai, Yoshua Bengio, Ian J Goodfellow, Erick Lavoie, Xavier Muller, Guillaume Desjardins, David Warde-Farley, et al. Unsupervised and transfer learning challenge: a deep learning approach. *Journal of Machine Learning Research-Proceedings Track*, 27:97–110, 2012.
- [81] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. In *Proceedings of International Conference of Learning Representations*, Scottsdale, Arizona, USA, 2013.
- [82] Tomas Mikolov, Martin Karafiát, Lukas Burget, Jan Cernocký, and Sanjeev Khudanpur. Recurrent neural network based language model. In *Proceedings of INTERSPEECH 2010*, 2010.
- [83] Tomas Mikolov, Quoc V Le, and Ilya Sutskever. Exploiting similarities among languages for machine translation. *CORR*, 2013.
- [84] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.
- [85] Dmitrijs Milajevs, Dimitri Kartsaklis, Mehrnoosh Sadrzadeh, and Matthew Purver. Evaluating neural word representations in tensor-based compositional settings. In *Proceedings of EMNLP*, 2014.
- [86] Jeff Mitchell and Mirella Lapata. Vector-based models of semantic composition. In *ACL*, pages 236–244, 2008.
- [87] Jeff Mitchell and Mirella Lapata. Composition in distributional models of semantics. *Cognitive science*, 34(8):1388–1429, 2010.
- [88] Andriy Mnih and Geoffrey E Hinton. A scalable hierarchical distributed language model. In *Advances in Neural Information Processing Systems*, pages 1081–1088, 2009.
- [89] Diego Mollá and José Luis Vicedo. Question answering in restricted domains: An overview. *Computational Linguistics*, 33(1):41–61, 2007.

- [90] Frederic Morin and Yoshua Bengio. Hierarchical probabilistic neural network language model. In *Proceedings of the international Workshop on Artificial Intelligence and Statistics*, pages 246–252, 2005.
- [91] Raymond H Myers. *Classical and Modern Regression with Applications*, volume 2. Duxbury Press Belmont, CA, 1990.
- [92] Douglas L Nelson, Cathy L McEvoy, and Thomas A Schreiber. The University of South Florida free association, rhyme, and word fragment norms. *Behavior Research Methods, Instruments, & Computers*, 36(3):402–407, 2004.
- [93] Donald A Norman. Memory, knowledge, and the answering of questions. *ERIC*, 1972.
- [94] Allan Paivio. Dual coding theory: Retrospect and current status. *Canadian Journal of Psychology*, 45(3):255, 1991.
- [95] Bo Pang and Lillian Lee. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings of the 42nd annual meeting on Association for Computational Linguistics*, page 271. Association for Computational Linguistics, 2004.
- [96] Bo Pang and Lillian Lee. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, pages 115–124. Association for Computational Linguistics, 2005.
- [97] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014)*, October 2014.
- [98] Nghia The Pham, Germán Kruszewski, Angeliki Lazaridou, and Marco Baroni. Jointly optimizing word representations for lexical and sentential tasks with the c-phrase model. In *Proceedings of ALC*, 2015.
- [99] Tamara Polajnar, Laura Rimell, and Stephen Clark. An exploration of discourse-based sentence spaces for compositional distributional semantics. In *Workshop on Linking Models of Lexical, Sentential and Discourse-level Semantics (LSD-Sem)*, page 1, 2015.

- [100] Stephen Roller and Sabine Schulte im Walde. A multimodal LDA model integrating textual, cognitive and visual modalities. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1146–1157, Seattle, Washington, USA, October 2013. Association for Computational Linguistics.
- [101] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, pages 1–42, 2014.
- [102] Iulian V Serban, Alessandro Sordoni, Yoshua Bengio, Aaron Courville, and Joelle Pineau. Building end-to-end dialogue systems using generative hierarchical neural network models. *arXiv preprint arXiv:1507.04808*, 2015.
- [103] Ryan Shaw, Anindya Datta, Debra VanderMeer, and Kaushik Dutta. Building a scalable database-driven reverse dictionary. *Knowledge and Data Engineering, IEEE Transactions on*, 25(3):528–540, 2013.
- [104] C Silberer and M Lapata. Learning grounded meaning representations with autoencoders. In *Proceedings of Association for Computational Linguistics*. ACL, 2014.
- [105] Carina Silberer, Vittorio Ferrari, and Mirella Lapata. Models of semantic representation with visual attributes. In *Proceedings of the 51th Annual Meeting of the Association for Computational Linguistics, Sofia, Bulgaria, August*, 2013.
- [106] Carina Silberer and Mirella Lapata. Grounded models of semantic representation. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 1423–1433. Association for Computational Linguistics, 2012.
- [107] Richard Socher, Jeffrey Pennington, Eric H Huang, Andrew Y Ng, and Christopher D Manning. Semi-supervised recursive autoencoders for predicting sentiment distributions. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 151–161. Association for Computational Linguistics, 2011.

- [108] Nitish Srivastava and Ruslan Salakhutdinov. Multimodal learning with deep boltzmann machines. In *NIPS*, pages 2231–2239, 2012.
- [109] Ilya Sutskever, James Martens, and Geoffrey E Hinton. Generating text with recurrent neural networks. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 1017–1024, 2011.
- [110] Ilya Sutskever, Oriol Vinyals, and Quoc VV Le. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112, 2014.
- [111] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. *arXiv preprint arXiv:1409.4842*, 2014.
- [112] Joseph Turian, Lev Ratinov, and Yoshua Bengio. 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, 2010.
- [113] Peter D Turney and Patrick Pantel. From frequency to meaning: Vector space models of semantics. *Journal of artificial intelligence research*, 37(1):141–188, 2010.
- [114] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103. ACM, 2008.
- [115] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, and Pierre-Antoine Manzagol. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *The Journal of Machine Learning Research*, 11:3371–3408, 2010.
- [116] Luis Von Ahn and Laura Dabbish. Labeling images with a computer game. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pages 319–326. ACM, 2004.

- [117] Ellen M Voorhees. Overview of the trec 2001 question answering track. *NIST special publication*, pages 42–51, 2002.
- [118] Ivan Vulić, Wim De Smet, and Marie-Francine Moens. Identifying word translations from comparable corpora using latent topic models. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies: short papers-Volume 2*, pages 479–484. Association for Computational Linguistics, 2011.
- [119] Ivan Vulic, Wim De Smet, and Marie-Francine Moens. Identifying word translations from comparable corpora using latent topic models. In *Proceedings of the Association for Computational Linguistics*, 2011.
- [120] Jason Weston, Antoine Bordes, Sumit Chopra, and Tomas Mikolov. Towards ai-complete question answering: a set of prerequisite toy tasks. *arXiv preprint arXiv:1502.05698*, 2015.
- [121] Janyce Wiebe, Theresa Wilson, and Claire Cardie. Annotating expressions of opinions and emotions in language. *Language resources and evaluation*, 39(2-3):165–210, 2005.
- [122] Katja Wiemer-Hastings and Xu Xu. Content differences for abstract and concrete concepts. *Cognitive Science*, 29(5):719–736, 2005.
- [123] Pengcheng Wu, Steven CH Hoi, Hao Xia, Peilin Zhao, Dayong Wang, and Chunyan Miao. Online multimodal deep similarity learning with application to image retrieval. In *Proceedings of the 21st ACM International Conference on Multimedia*, pages 153–162. ACM, 2013.
- [124] Matthew D. Zeiler. Adadelta: An adaptive learning rate method. In *arXiv preprint arXiv:1212.5701*, 2012.
- [125] Michael Zock and Slaven Bilac. Word lookup on the basis of associations: From an idea to a roadmap. In *Proceedings of the ACL Workshop on Enhancing and Using Electronic Dictionaries*, 2004.



# **Appendix A**

## **Extra Information**

Some more text ...

