# A Comparative Study of Word Embeddings
# for Reading Comprehension

**Bhuwan Dhingra, Hanxiao Liu, Ruslan Salakhutdinov, William W. Cohen**
School of Computer Science
Carnegie Mellon University, Pittsburgh, USA
{bdhingra, hanxiaol, rsalakhu, wcohen}@cs.cmu.edu

## Abstract

The focus of past machine learning research for Reading Comprehension tasks has been primarily on the design of novel deep learning architectures. Here we show that seemingly minor choices made on (1) the use of pre-trained word embeddings, and (2) the representation of out-of-vocabulary tokens at test time, can turn out to have a larger impact than architectural choices on the final performance. We systematically explore several options for these choices, and provide recommendations to researchers working in this area.

## 1 Introduction

Systems that can read documents and answer questions about their content are a key language technology. The field, which has been termed *Reading Comprehension* (RC), has attracted a tremendous amount of interest in the last two years, primarily due to the introduction of large-scale annotated datasets, such as CNN (Hermann et al., 2015) and SQuAD (Rajpurkar et al., 2016).

Powerful statistical models, including deep learning models (also termed as *readers*), have been proposed for RC, most of which employ the following recipe: (1) Tokens in the document and question are represented using word vectors obtained from a lookup table (either initialized randomly, or from a pre-trained source such as *GloVe* (Pennington et al., 2014)). (2) A sequence model such as LSTM (Hochreiter and Schmidhuber, 1997), augmented with an attention mechanism (Bahdanau et al., 2014), updates these vectors to produce contextual representations. (3) An output layer uses these contextual representations to locate the answer in the document. The focus so far in the literature has been on steps (2) and
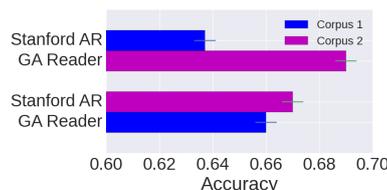


Figure 1: Test set accuracies and std error on the Who-Did-What dataset for Stanford AR and GA Reader, trained after initializing with word vectors induced from different corpora. Without controlling for the initialization method, different conclusions may be drawn about which architecture is superior. Corpus 1: BookTest dataset (Bajgar et al., 2016), Corpus 2: Wikipedia + Gigaword.

(3), and several novel architectures have been proposed (see Section 2.1).

In this work, we show that seemingly minor choices made in step (1), such as the use of pre-trained word embeddings and the handling of out-of-vocabulary tokens at test time, can lead to substantial differences in the final performance of the reader. These differences are usually much larger than the gains reported due to architectural improvements. As a concrete example, in Figure 1 we compare the performance of two RC models—Stanford Attentive Reader (AR) (Chen et al., 2016) and Gated Attention (GA) Reader (Dhingra et al., 2016)—on the Who-Did-What dataset (Onishi et al., 2016), initialized with word embeddings trained on different corpora. Clearly, comparison between architectures is meaningful only under a controlled initialization method.

To justify our claims, we conduct a comprehensive set of experiments comparing the effect of utilizing embeddings pre-trained on several corpora. We experiment with RC datasets from different domains using different architectures, and obtain consistent results across all settings. Based on our findings, we recommend the use of certain pre-trained GloVe vectors for initialization. These consistently outperform other off-the-shelf

embeddings such as *word2vec* [1] (Mikolov et al., 2013), as well as those pre-trained on the target corpus itself and, perhaps surprisingly, those trained on a large corpus from the same domain as the target dataset.

Another important design choice is the handling of out-of-vocabulary (OOV) tokens at test time. A common approach (e.g. (Chen et al., 2016; Shen et al., 2016)) is to replace infrequent words during training with a special token UNK, and use this token to model the OOV words at the test phase. In reading comprehension, where the target answers are often rare words, we find that this approach leads to significantly worse performance in certain cases. A superior strategy is to assign each OOV token either a pre-trained, if available, or a random but unique vector at test time. We discuss and compare these two strategies, as well a mixed variant of the two, in Section 3.2.

## 2 Background

### 2.1 RC Datasets & Models

Many datasets aimed at measuring the performance of RC have been proposed (Nguyen et al., 2016; Trischler et al., 2016). For our purposes, we pick two of these benchmarks from different domains – Who-Did-What (WDW) (Onishi et al., 2016) constructed from news stories, and the Children's Book Test (CBT) (Hill et al., 2015) constructed from children's books. For CBT we only consider the questions where the answer is a named entity (CBT-NE). Several RC models based on deep learning have been proposed (Cui et al., 2016; Munkhdalai and Yu, 2016; Sordoni et al., 2016; Shen et al., 2016; Kobayashi et al., 2016; Henaff et al., 2016; Wang and Jiang, 2016; Wang et al., 2016; Seo et al., 2016; Xiong et al., 2016; Yu et al., 2016). For our experiments we pick two of these models: the simple, but competitive, Stanford AR, and the high-performing GA Reader.

**Stanford AR:** The Stanford AR consists of single-layer Bidirectional GRU encoders for both the document and the query, followed by a bilinear attention operator for computing a weighted average representation of the document. The original model, which was developed for the anonymized CNN / Daily Mail datasets, used an output lookup table $W_a$ to select the answer. However, without anonymization the number of answer candi-

| Emb. | Corpus | Domain | Size | Vocab |
|------|--------|--------|------|-------|
| OTS | Wiki + Gigaword / GoogleNews | Wiki / News | 6B / 100B | 400K / 3M |
| WDW | Who-Did-What | News | 50M | 91K |
| BT | BookTest | Fiction | 8B | 1.2M |
| CBT | Children's BookTest | Fiction | 50M | 48K |

Table 1: Details of corpora used for training word embeddings. **OTS:** Off-The-Shelf embeddings provided with GloVe / word2vec. Corpus size is in # of tokens.

dates can become very large. Hence, we instead select the answer from the document representation itself, followed by an *attention sum* mechanism (Kadlec et al., 2016). This procedure is very similar to the one used in GA Reader, and is described in detail in Appendix A.

**GA Reader:** The GA Reader is a multi-hop architecture which updates the representation of document tokens through multiple bidirectional GRU layers (Cho et al., 2014). At the output of each intermediate layer, the token representations are re-weighted by taking their element-wise product with an attention-weighted representation of the query. The outputs of the final layer are further matched with the query representation with an inner product to produce a distribution over the candidate answers in the document, and multiple mentions are aggregated using attention sum. We use the publicly available code[2] with the default hyperparameter settings of (Dhingra et al., 2016), detailed in Appendix B.

### 2.2 Word Embedding methods

The two most popular methods for inducing word embeddings from text corpora are *GloVe* (Pennington et al., 2014) and *word2vec* (Mikolov et al., 2013). These packages also provide off-the-shelf (OTS) embeddings trained on large corpora[3]. While the GloVe package provides embeddings with varying sizes (50-300), word2vec only provides embeddings of size 300. This is an important difference, which we discuss in detail later. We also train three additional embeddings, listed in Table 1, including those trained on the target datasets themselves. In summary, we test with two in-domain corpora for WDW: one large (OTS) and one small (WDW), and two in-domain corpora for CBT: one large (BT) and one small (CBT).

---

[1] https://code.google.com/archive/p/word2vec/

[2] https://github.com/bdhingra/ga-reader

[3] The word2vec package contains embeddings for both capitalized and lowercase words. We convert all words to lowercase, and if a word has both lowercase and uppercase embeddings we use the lowercase version.

When training embeddings, we set hyperparameters to their default values in the provided packages (see Appendix B for details). This is by no means an optimal choice, in fact previous studies (Levy et al., 2015) have shown that hyperparameter choices may have a significant impact on downstream performance. However, training a single RC model can take anywhere from several hours to several days, and tuning hyperparameters for the embedding method on this downstream task is both infeasible and rarely done in practice. Instead, our objective is to provide guidelines to researchers using these methods out-of-the-box.

## 3 Experiments and Results

### 3.1 Comparison of Word Embeddings

We repeat each experiment twice with different random seeds and report the average test set accuracy across the two runs. Figure 2 shows a comparison of the RC performance for GA Reader and Stanford AR after initializing with various pre-trained embeddings, and also after initializing randomly. We see consistent results across the two datasets and and the two models.

The first observation is that using embeddings trained on the right corpora can improve anywhere from 3-6% over random initialization. However, the corpus and method used for pre-training are important choices: for example word2vec embeddings trained on CBT perform worse than random. Also note that in every single case, GloVe embeddings outperform word2vec embeddings trained on the same corpora. It is difficult to claim that one method is better than the other, since previous studies (Levy et al., 2015) have shown that these methods are sensitive to hyperparameter tuning. However, if used out-of-the-box, GloVe seems to be the preferred method for pre-training.

The single best performance is given by off-the-shelf GloVe embeddings ($d = 100$) in each case, which outperform off-the-shelf word2vec embeddings ($d = 300$). To understand if the difference comes from the differing dimension sizes, we plot the performance of GloVe embeddings as the dimension size is increased in Figure 3 (left). Performance drops as the embedding dimension size is increased (most likely due to over-fitting); however even at $d = 300$, GloVe embeddings outperform word2vec embeddings.

On both test datasets embeddings trained on formal domains, like news (OTS, WDW), perform at least as well as those trained on informal ones, like fiction (BT, CBT). This is surprising for CBT-NE dataset which is itself constructed from the informal domain of children's books. For example, WDW (50M tokens) does significantly better than CBT-NE (50M tokens) in 3 out of the 4 cases, and also significantly better than the much larger BT (8B tokens) in one setting (and comparably in other settings). A key distinguishing feature between these two domains is the fraction of text composed of stopwords [4] – WDW consists of 54% stopwords while BT consists of 68% stopwords. Both GloVe and word2vec induce word vectors by minimizing the Euclidean distance between vectors of frequently co-occurring words. Co-occurrence with stopwords, however, provides little meaningful information about the semantics of a particular word, and hence corpora with a high percentage of these may not produce high-quality vectors. This effect may be mitigated during pre-training by either, (1) removing a fraction of stopwords from the corpus, (2) increasing the window size for counting co-occuring words. Figure 3 (right) shows the effect of both these methods on downstream RC performance. There is an improvement as the fraction of stopwords decreases or the window size increases, upto a certain limit. In fact, with proper tuning, BT embeddings give roughly the same performance as OTS GloVe, emphasizing the importance of hyperparameter tuning when training word vectors. There is also evidence that stopword removal can be beneficial when training word vectors, however this needs further verification on other downstream tasks.

### 3.2 Handling OOV tokens

In this section we study some common techniques for dealing with OOV tokens at test time. Based on the results from the previous section, we conduct this study using only the off-the-shelf GloVe pre-trained embeddings. Let the training, test and GloVe vocabularies be denoted by $V^T$, $V^E$ and $V^G$ respectively. Also define $V_n^T = \{t \in V^T : \#t > n\}$ where $\#t$ denotes the count of token $t$ in the training corpus. Before training a neural network for RC, the developer must first decide on the set of words $V$ which will be assigned word vectors. Any token outside $V$ is treated as an OOV token (denoted by UNK) and is assigned the same

---

[4] We use the list of stopwords available at http://research.microsoft.com/en-us/um/redmond/projects/mctest/data/stopwords.txt
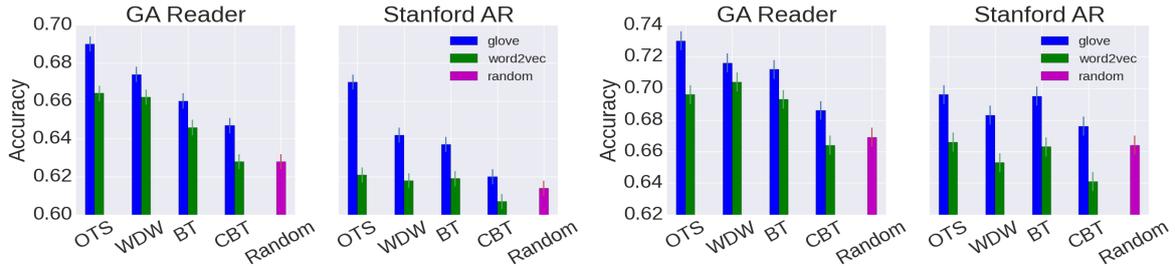
**Figure 2:** Test set accuracy and std error for GA Reader and Stanford AR on WDW (left, middle-left) and CBT-NE (middle-right, right) when trained after initializing with pre-trained embeddings induced from different corpora (Table 1), or randomly
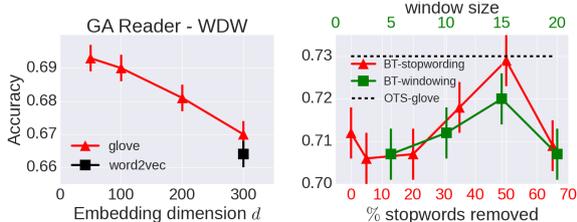


**Figure 3:** Test set accuracy and std error on **left:** WDW when initialized with off-the-shelf GloVe embeddings of different sizes, **right:** CBT-NE when initialized with embeddings trained on BT corpus after removing a fraction of stopwords (red), or using different window sizes (green).
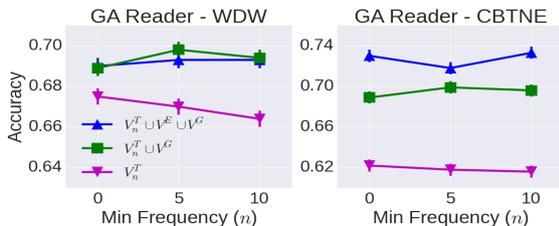


**Figure 4:** Test set accuracy and std error for GA Reader on WDW (**left**) and CBT-NE (**right**) when trained after assigning different cuts of the vocabulary with word vectors. *Min frequency* refers to the minimum count of a word type for it to be included in the vocabulary.

fixed vector.

By far the most common technique in NLP literature (e.g. (Chen et al., 2016; Shen et al., 2016)) for constructing this vocabulary is to decide on a minimum frequency threshold $n$ (typically 5-10) and set $V = V_n^T$. Out of these, vectors for those which also appear in $V^G$ are initialized to their GloVe embeddings, and the rest are randomly initialized. Remaining tokens in $V^T$ and those in $V^E - V^T$ are all assigned the UNK vector, which is itself updated during training. This method ignores the fact that many of the words assigned as UNK may have already trained embeddings available in $V^G$. Hence, here we propose another strategy of constructing the vocabulary as $V = V_n^T \cup V^G$. Then at test time, any new token would be assigned its GloVe vector if it exists, or the vector for UNK. A third approach, used in (Dhingra et al., 2016), is motivated by the fact that many of the RC models rely on computing fine-grained similarity between document and query tokens. Hence, instead of assigning all OOV tokens a common UNK vector, it might be better to assign them untrained but unique random vectors. This can be done by setting the vocabulary to $V = V_n^T \cup V^E \cup V^G$. Hence, at test time any new token will be assigned its GloVe vector if it exists, or a random vector. Note that for this approach access to $V^E$ at training time is not needed.

Figure 4 shows a comparison of all three approaches with varying $n$ for the GA Reader on WDW and CBT-NE datasets. A gap of $3\%$ and $11\%$ between the best and worst setting for WDW and CBT-NE respectively clearly indicates the importance of using the correct setting. The commonly used method of setting $V = V_n^T$ is not a good choice for RC, and gets worse as $n$ is increased. It performs particularly poorly for the CBT-NE dataset, where $\sim 20\%$ of the test set answers do not appear in the training set (compared to only $\sim 1.5\%$ in WDW). The other two approaches perform comparably for WDW, but for CBT-NE assigning random vectors rather than UNK to OOV tokens gives better performance. This is also easily explained by looking at fraction of test set answers which do not occur in $V_0^T \cup V^G$ – it is $\sim 10\%$ for CBT-NE, and $< 1\%$ for WDW. Since in general it is not possible to compute these fractions without access to the test set, we recommend setting $V = V_n^T \cup V^E \cup V^G$.

## 4 Conclusions

We have shown that the choice of pre-trained embeddings for initializing word vectors has a significant impact on the performance of neural models for reading comprehension. So does the method for handling OOV tokens at test time. We argue that different architectures can only be compared when these choices are controlled for. Based on our experiments, we recommend the use of off-the-shelf GloVe embeddings, and assigning pre-trained GloVe vectors, if available, or random but unique vectors to OOV tokens at test time.

## References

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* .

Ondrej Bajgar, Rudolf Kadlec, and Jan Kleindienst. 2016. Embracing data abundance: Booktest dataset for reading comprehension. *arXiv preprint arXiv:1610.00956* .

Danqi Chen, Jason Bolton, and Christopher D Manning. 2016. A thorough examination of the cnn/daily mail reading comprehension task. *arXiv preprint arXiv:1606.02858* .

Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078* .

Yiming Cui, Zhipeng Chen, Si Wei, Shijin Wang, Ting Liu, and Guoping Hu. 2016. Attention-over-attention neural networks for reading comprehension. *arXiv preprint arXiv:1607.04423* .

Bhuwan Dhingra, Hanxiao Liu, William W Cohen, and Ruslan Salakhutdinov. 2016. Gated-attention readers for text comprehension. *arXiv preprint arXiv:1606.01549* .

Mikael Henaff, Jason Weston, Arthur Szlam, Antoine Bordes, and Yann LeCun. 2016. Tracking the world state with recurrent entity networks. *arXiv preprint arXiv:1612.03969* .

Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *Advances in Neural Information Processing Systems*. pages 1684–1692.

Felix Hill, Antoine Bordes, Sumit Chopra, and Jason Weston. 2015. The goldilocks principle: Reading children's books with explicit memory representations. *arXiv preprint arXiv:1511.02301* .

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9(8):1735–1780.

Rudolf Kadlec, Martin Schmid, Ondrej Bajgar, and Jan Kleindienst. 2016. Text understanding with the attention sum reader network. *arXiv preprint arXiv:1603.01547* .

Diederik Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* .

Sosuke Kobayashi, Ran Tian, Naoaki Okazaki, and Kentaro Inui. 2016. Dynamic entity representations with max-pooling improves machine reading. In *NAACL-HLT* .

Omer Levy, Yoav Goldberg, and Ido Dagan. 2015. Improving distributional similarity with lessons learned from word embeddings. *Transactions of the Association for Computational Linguistics* 3:211–225.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*. pages 3111–3119.

Tsendsuren Munkhdalai and Hong Yu. 2016. Neural semantic encoders. *arXiv preprint arXiv:1607.04315* .

Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. Ms marco: A human generated machine reading comprehension dataset. *arXiv preprint arXiv:1611.09268* .

Takeshi Onishi, Hai Wang, Mohit Bansal, Kevin Gimpel, and David McAllester. 2016. Who did what: A large-scale person-centered cloze dataset. *EMNLP* .

Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*. pages 1532–1543. http://www.aclweb.org/anthology/D14-1162.

Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250* .

Minjoon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. 2016. Bidirectional attention flow for machine comprehension. *arXiv preprint arXiv:1611.01603* .

Yelong Shen, Po-Sen Huang, Jianfeng Gao, and Weizhu Chen. 2016. Reasonet: Learning to stop reading in machine comprehension. *arXiv preprint arXiv:1609.05284* .

Alessandro Sordoni, Phillip Bachman, and Yoshua Bengio. 2016. Iterative alternating neural attention for machine reading. *arXiv preprint arXiv:1606.02245* .

Adam Trischler, Tong Wang, Xingdi Yuan, Justin Harris, Alessandro Sordoni, Philip Bachman, and Kaheer Suleman. 2016. Newsqa: A machine comprehension dataset. *arXiv preprint arXiv:1611.09830* .

Shuohang Wang and Jing Jiang. 2016. Machine comprehension using match-lstm and answer pointer. *arXiv preprint arXiv:1608.07905* .

Zhiguo Wang, Haitao Mi, Wael Hamza, and Radu Florian. 2016. Multi-perspective context matching for machine comprehension. *arXiv preprint arXiv:1612.04211* .

Caiming Xiong, Victor Zhong, and Richard Socher. 2016. Dynamic coattention networks for question answering. *arXiv preprint arXiv:1611.01604* .

Yang Yu, Wei Zhang, Kazi Hasan, Mo Yu, Bing Xiang, and Bowen Zhou. 2016. End-to-end answer chunk extraction and ranking for reading comprehension. *arXiv preprint arXiv:1610.09996* .

## A    Answer Selection for Stanford AR

Using the notation from (Chen et al., 2016), let $\tilde{p}_1, \tilde{p}_2, \ldots, \tilde{p}_m$ be the contextual embeddings of the tokens in the document, and let $o$ be the attention-weighted document representation, then we compute the probability that token $i$ answers the question as:

$$P(a = d_i | d, q) = s_i = \text{softmax}(\tilde{p}_i^T o) \quad (1)$$

The probability of a particular candidate $c \in \mathcal{C}$ as being the answer is then computed by aggregating the probabilities of all document tokens which appear in $c$ and renormalizing over the candidates:

$$\Pr(c|d, q) \propto \sum_{i \in \mathbb{I}(c,d)} s_i \quad (2)$$

where $\mathbb{I}(c, d)$ is the set of positions where a token in $c$ appears in the document $d$.

## B    Hyperparameter Details

For the WDW dataset we use hidden state size $d = 128$ for the GRU and dropout with $p = 0.3$. For CBT-NE dataset we use $d = 128$ and dropout with $p = 0.4$. The Stanford AR has only 1 layer as proposed in the original paper, while the GA Reader has 3 layers. For Stanford AR dropout is applied to the input of the layer, and for GA Reader it is applied in between layers. Embeddings sizes for the word vectors were set to $d_w = 100$ for all experiments, except those using off-the-shelf word2vec embeddings. To enable a fair comparison, we utilize the *qe-comm* feature for Stanford AR, which was used in the implementation of GA Reader. Since our purpose is to study the effect of word vectors, we do not use character embeddings in our experiments.

We train the models using the ADAM (Kingma and Ba, 2014) optimizer with an initial learning rate of 0.0005, which is halved every epoch after the first 3 epochs. We track performance on the validation set, and select the model with the highest validation accuracy for testing.

When training word vectors we retain the default settings provided with the GloVe and word2vec packages, with the only exception that window size was set to 10 for both (to ensure consistency). For word2vec, we used skip-gram architecture with hierarchical softmax, and subsampled frequent words with a threshold $10^{-3}$ (see (Mikolov et al., 2013) for details). For GloVe, we used 15 iterations when training on the small WDW and CBT corpora, and 50 iterations for the large BT corpus. In any corpus, words occurring less than 5 times were filtered before training the word vectors.