
Exploring Efficacy of Embeddings on Relation Network for Natural Language Question Answering Task

Anonymous Author(s)

Affiliation

Address

email

1 Introduction

Deep learning has made it possible to do classification of objects in images and translation of languages, often with incredible accuracy. This is achieved due to the ability of neural networks to pick out important patterns that are inconceivable to the human eye, from large quantities of labeled data. However, just being able to learn patterns is not sufficient as it is not the only ability associated to intelligence; reasoning is another essential ability [1] that separates humans from machines. Hence, in recent years there is much work on reasoning related research, like visual reasoning [4, 9] where the machine is able to give an answer given an image and a visual question about the image, and text-based question answering [9] where the machine is able to answer a question based on the earlier sentences given to it.

In this project, we focus on the text-based question answering task using relation network (RN) [9] on the bAbI dataset [10]. RNs are networks that are designed based on relational reasoning, where its capacity to compute relations is baked into the architecture without having the neural network to learn it.

For any neural approach for natural language processing, word and sentence embeddings are indispensable. They allow us to represent words and sentences whose original forms are strings, as vectors which then we can feed it into an artificial neural network. There are various ways to embed words, and while unsupervised representations have been the more commonly used approach, using the assumption that you can tell a word by the company it keeps, there is an increased focus on supervised representations and also multi-task learning of representations. In our project, we explore how different types of embeddings, in particular, how the traditional unsupervised representations compare up to representations obtained from multi-task learning.

Our experiments involve using two different representations, the first one being the approach used by the original RN paper [9], to embed the context and questions into sentence embeddings using LSTMs, and the other uses the universal sentence encoder (USE) [2] to embed the context and questions into sentence embeddings. In the RN paper, sentence embeddings are called objects which the RN is uses to learn the relation between them. We then compare the performance of different embeddings on RN for bAbI question answering task.

2 Task

2.1 bAbI

The bAbI dataset is a pure text-based question answering (QA) dataset that contains a total of 20 tasks. Each task corresponds to a particular type of reasoning, such as deduction, induction and counting. Every question is associated with a set of supporting facts, which provides the context for

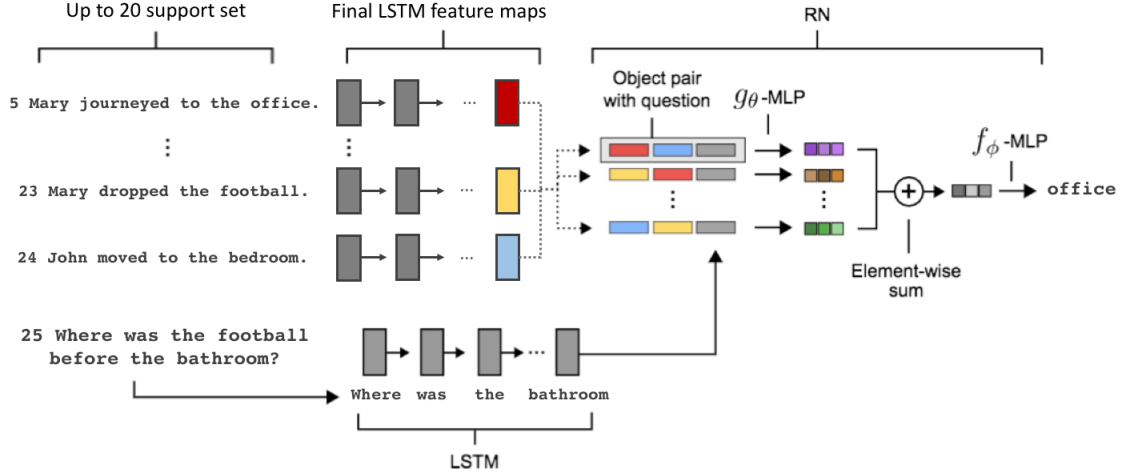


Figure 1: Text based QA architecture. Contexts and questions are processed with an LSTM to produce a set of context and question embedding. Objects, colored yellow, red, and blue, are constructed using LSTMs or USE. The RN considers relations across all pairs of objects, conditioned on the question embedding, and integrates all these relations to answer the question. Our alternative approach substitutes the LSTM with the USE to produce the embeddings.

the question being asked. An example “Sandra picked up the football” and “Sandra went to the office” support the question “Where is the football?”, which we humans can arrive at an easy at the answer “office”. A task is considered to be successfully passed if it attains an accuracy of 95% or higher.

2.2 Two Supporting Fact Task

Task 2 of the bAbI tasks requires chaining of two or three supporting facts to answer the question. As such, to answer the question “Where is the football?”, it has to possibly be able to link information from the sentences “Mary moved to the bathroom”, “Mary picked up the football there” and “Mary went back to the garden” for it to conclude that the football is at the garden. This tasks makes it challenging for the neural network as it requires some form of memory for it to be able to link previously acquired knowledge to answer the question.

The RN succeeded on the basic induction task, which proved difficult for Sparse DNC, Differentiable Neural Computer (DNC) [3] and Sparse DNC [8]. However, it missed the 95% mark for the “two supporting fact” and “three supporting fact” tasks, which architectures with memory like Memory Networks [11], DNC and sparse DNC excelled. Thus we focused on task 2 for our project, to investigate if a better representation, in a form the USE obtained from multi-task learning, fare better than unsupervised representations.

3 Model

1. Overview of the original RN model, comment on the strength and weaknesses
2. Modifications to the RN model that will help improve the accuracy of the task. Motivations for the modifications.
3. (Optional) A paragraph on USE?
4. How long we take to train our model and the train/test accuracy, loss values etc. Use original RN paper as a guideline of what numbers to show.

3.1 Relation Network

The RN is a neural network module which is designed to do relational reasoning. It is a composite function whose form is given by:

$$RN(O) = f_{\phi} \left(\sum_{i,j} g_{\theta}(o_i, o_j, q) \right) \quad (1)$$

where the input is a set of “objects” $O = \{o_1, \dots, o_n\}$, $o_i \in \mathbb{R}^m$ is the i^{th} object, and f_{ϕ} and g_{θ} are function with learnable parameters ϕ and θ respectively. f_{ϕ} and g_{θ} are multi-layer perceptrons (MLP), where g_{θ} learns the relation between any two given objects and f_{ϕ} maps the relations to one of the many possible answers. In the next part, we will discuss how using different embeddings lead to different implementations of our model.

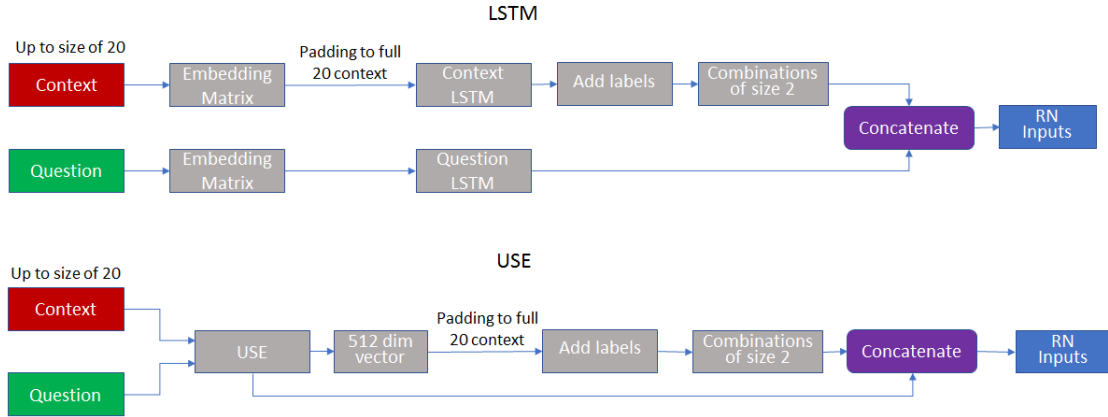


Figure 2: Embedding Pipelines. Padding are done at different stages of the pipeline for the different embeddings in which padding is done first for the LSTM, whereas padding is done after passing through the USE.

3.1.1 LSTM Embeddings

The bAbI suite of tasks requires transformation of the natural language inputs into a set of objects. In the original RN paper, to provide the prior knowledge for the RN to base on, 20 sentences in the support set that were prior to the question of interest were used as the context for answering the question, and each sentence was processed word-by-word with an 32 unit LSTM to obtain a sentence embedding, after which it is tagged with labels to indicate their relative position in the support set to obtain the objects that are the inputs to our RN. A separate LSTM with 32 units was used to process the question. Lastly, the softmax output was optimized with a cross-entropy loss function using the Adam optimizer with a learning rate of $2e^{-4}$.

From (1) we see that for a context of size k , there are $k(k-1)/2$ unique pairings of the objects. In its implementation, for questions with context sizes that are smaller than 20, we pad it with zeros such that we have a full context size of 20, then let the LSTM process the words to get the context embeddings. The labels that tagged the relative position in the support set, which are just one-hot vectors indicating the position of the sentence in the support set, are then concatenated to the embeddings. Lastly, for each given question embedding, it is concatenated to all possible combination of size 2 to obtain the inputs to the RN, which has dimensions (190, 136), the first value reflecting the number of possible combinations and the other the sum of the dimension of the objects and question embeddings.

3.1.2 USE Embeddings

For the USE embeddings, we followed the original RN paper of using up to 20 sentences in the support set as the prior for answering the question. The string form of the raw context and question

86 sentences are fed into the USE which outputs a 512 dimensional vector sentence embedding for each
87 sentence. As such we see that the USE embeddings gives the contexts and questions a much richer
88 representation as compared to the LSTM embeddings.

89 One main difference from the LSTM embedding approach is that we do not pad the context to the full
90 size of 20 before passing through the USE, compared to padding the context to 20 sentences before
91 letting the LSTM process the words. We note that to do something similar in the USE context, we
92 have to pad the sentences to the same length with a special character, for example “null”. Instead
93 we keep the sentence length as it is, pass it through the USE to obtain the sentence embeddings and
94 for context with size less than 20, pad it sufficiently with zeros we get 20 contexts. The addition of
95 the one-hot vector labels, we only concatenate the one-hot vectors for those non-trivial contexts, for
96 those padded contexts, the labels are just zero vectors of dimension 20. From experiments, if we
97 were to label the trivial contexts with the one-hot vectors, it gives a very poor performance. As the
98 dimensions of the representation of each sentence for both context and questions is increased from 32
99 to 512, our inputs to the RN has dimensions (190, 1536).

100 3.2 Embeddings

101 (Short paragraph on embeddings, bert, Elmo, glove, etc.)

102 There are various ways to embed words, from common unsupervised models like word2vec [5] and
103 GloVe [6]. Early this year we have ELMo [7], which improved the state of the art embeddings. ELMo
104 differs from other unsupervised embeddings by having their inputs to be characters instead of words,
105 allowing them to take advantage of sub-word units to compute meaningful representations even for
106 out-of-vocabulary words. We shall now discuss the two different embeddings we have used in our
107 experiments whose efficacy we are investigating.

108 3.2.1 LSTM

109 For the LSTM embeddings, a randomly initialized word embedding matrix is used to assign each
110 word to a vector representation, which is then passed through a 32 unit LSTM
111 hashing on the word level, the lstm to obtain an embedding on the sentence level.

112 3.2.2 Universal Sentence Encoder

113 paper, tensorflow blog

114 4 Results

115 We have made several comparison between different types of embeddings, different sample sizes,
116 different types of inputs to feed to the RN module and different modifications to the RN module
117 with our implementation of USE embeddings. We compared the accuracy of the test set for different
118 modifications to the RN module and compared the accuracy of the validation set for the rest of the
119 comparisons.

120 For the experiments, all models were run with batch size of 64, with 200 epochs and validation set
121 was taken from the last 10% of the training set. For the original RN module, $f_\phi \left(\sum_{i,j} g_\theta(o_i, o_j, q) \right)$,
122 g_θ is a four-layer MLP 256 units per layer and f_ϕ is a three-layer MLP consisting of 256, 512 and
123 6 units, where the final softmax output was optimized with a cross-entropy loss function using the
124 Adam optimizer with a learning rate of $2e^{-4}$. All models are run on original RN module except for
125 the different modifications to the RN module. The results and weights used are all based on the epoch
126 that gave the best validation accuracy.

127 4.1 Comparison between Different Embeddings

128 We have made comparisons between the two different embeddings of the sentences of the contexts,
129 questions and answers from all data from the training set of Task 2:

- 130 • Model LSTM: Use LSTM to embed the sentences of the contexts, questions and answers,
131 and then to pass to the RN module

- Model USE: Use USE to embed the sentences of the contexts, questions and answers, and then to pass to the RN module

Model	Training Accuracy (%)	Training Accuracy (%)
LSTM		
USE		

From the results, we can see that based on training solely on the training set of Task 2, using USE embedding provides a much better representation than LSTM embedding in training our RN module. Thus, we will be running the rest of the models with USE embedding.

4.2 Comparison between Different Sample Sizes

We ran the original model and compared between 2 different sample sizes of the training set:

- First 1000 contexts of the training set
- Full 9995 contexts of the training set

Model	Training Accuracy (%)	Training Accuracy (%)
First 1000	96.2	35.0
Full 9995	100	87.1

Based on the validation accuracy, we can see that it is necessary to include all data in the training set for training of the models. Thus, we will be running the rest of the models based on the full 9995 contexts.

4.3 Comparison between Different Modification of the RN Module

For LSTM embedding, each sentence and question had been embedded to a 32-unit object, with sentences from the contexts and questions being processed through separate LSTMs. For each question embedding, it is paired with a combination of 2 sentence embeddings from up to 20 prior sentences, with indexes. Batches of 190 (20 sentences choose 2) combinations of dimension 136 vectors (32-unit sentence i embedding + one hot 20-dimension vector for index of sentence i + 32-unit sentence j embedding + one hot 20-dimension vector for index of sentence j + 32-unit question embedding) with batch size of 64 were passed through the RN module.

On the other hand, each sentence and question had been embedded to a 512-unit object with the pre-trained USE. As a result, instead of having 190 combinations of dimension 136 vectors with batch size of 64 being passed through the RN module, we would have combinations of dimension 1576 (512-unit sentence i embedding + one hot 20-dimension vector for index of sentence i + 512-unit sentence j embedding + one hot 20-dimension vector for index of sentence j + 512-unit question embedding). From this we can see that there is an increase in the dimension of each input, and hence, the complexity of each input before passing through the RN module. Thus, we believed by increasing the number of layers or/and number of neurons in each layer would improve the accuracy of the test set as the RN module would have more components to learn from.

We compared between different modifications of the RN model and checked for the accuracy of the test set:

- Model ori: Original RN module ($g : [256 \times 4], f : [256, 512, 6]$)
- Model g6: Increase number of layers of g from 4 to 6 ($g : [256 \times 6], f : [256, 512, 6]$)
- Model g8: Increase number of layers of g from 4 to 8 ($g : [256 \times 8], f : [256, 512, 6]$)
- Model 512: Increase number of neurons in each layer of g from 256 to 512 ($g : [512 \times 4], f : [256, 512, 6]$)
- Model 512f512: Increase number of neurons in each layer of g from 256 to 512 and first layer of f from 256 to 512 ($g : [512 \times 4], f : [512, 512, 6]$)
- Model g6x512: Increase number of layers of g from 4 to 6 and increase number of neurons in each layer of g from 256 to 512 ($g : [512 \times 6], f : [256, 512, 6]$)

- Model g6x512f512: Increase number of layers of g from 4 to 6, increase number of neurons in each layer of g from 256 to 512 and first layer of f from 256 to 512 ($g : [512 \times 6]$, $f : [512, 512, 6]$)
- Model double: Increase number of layers of g from 4 to 6 and double the number of neurons in all layers except for the last layer of f ($g : [512 \times 6]$, $f : [512, 1024, 6]$)

Model	Training Accuracy (%)	Training Accuracy (%)	
ori	100	87.1	85.5
g6	100	89.1	88.4
g8	18.4	22.2	20.0
512	100	88.2	85.0
512f512	99.9	87.9	88.3
g6x512	100	87.8	88.5
g6x512f512	99.7	88.3	89.0
double	100	88.2	88.4

Based on the test accuracies of the above models, we can see that increasing the number of layers of g function from 4 to 6 increased the test accuracy but drastically decreased the accuracy when it was increased to 8 layers. There is also improvement in the test accuracy when the number of neurons for each layer of g function from 256 to 512. The above table shows that the best model for our USE implementation is Model g6x512f512, which is to increase both the number of layers from 4 to 6 and the number of neurons from 256 to 512 for each layer of the g function, together with an increase of number of neurons from 256 to 512 for the first layer of the f function.

5 Discussion and Conclusions

References

- [1] L. Bottou. From Machine Learning to Machine Reasoning. *Arxiv preprint arXiv11021808*, page 15, 2011. ISSN 0885-6125. doi: 10.1007/s10994-013-5335-x. URL <http://arxiv.org/abs/1102.1808>.
- [2] D. Cer, Y. Yang, S.-y. Kong, N. Hua, N. Limtiaco, R. S. John, N. Constant, M. Guajardo-Cespedes, S. Yuan, C. Tar, Y.-H. Sung, B. Strope, and R. Kurzweil. Universal Sentence Encoder. 2018. ISSN 1042-1629. doi: 10.1007/s11423-014-9358-1. URL <http://arxiv.org/abs/1803.11175>.
- [3] A. Graves, G. Wayne, M. Reynolds, T. Harley, I. Danihelka, A. Grabska-Barwińska, S. G. Colmenarejo, E. Grefenstette, T. Ramalho, J. Agapiou, A. P. Badia, K. M. Hermann, Y. Zwols, G. Ostrovski, A. Cain, H. King, C. Summerfield, P. Blunsom, K. Kavukcuoglu, and D. Hassabis. Hybrid computing using a neural network with dynamic external memory. *Nature*, 538(7626): 471–476, 2016. ISSN 14764687. doi: 10.1038/nature20101. URL <http://dx.doi.org/10.1038/nature20101>.
- [4] J. Johnson, B. Hariharan, L. V. D. Maaten, J. Hoffman, L. Fei-Fei, C. L. Zitnick, and R. Girshick. Inferring and Executing Programs for Visual Reasoning. *Proceedings of the IEEE International Conference on Computer Vision*, 2017-Octob:3008–3017, 2017. ISSN 15505499. doi: 10.1109/ICCV.2017.325.
- [5] T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781, 2013. URL <http://dblp.uni-trier.de/db/journals/corr/corr1301.html#abs-1301-3781>.
- [6] J. Pennington, R. Socher, and C. D. Manning. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014. URL <http://www.aclweb.org/anthology/D14-1162>.
- [7] M. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*,

- 215 *Volume 1 (Long Papers)*, pages 2227–2237. Association for Computational Linguistics, 2018.
 216 doi: 10.18653/v1/N18-1202. URL <http://aclweb.org/anthology/N18-1202>.
- 217 [8] J. W. Rae, J. J. Hunt, T. Harley, I. Danihelka, A. Senior, G. Wayne, A. Graves, and T. P. Lillicrap.
 218 Scaling Memory-Augmented Neural Networks with Sparse Reads and Writes. (Nips), 2016.
 219 ISSN 10495258. doi: 10.1016/j.molimm.2007.07.004. URL [http://arxiv.org/abs/1610.](http://arxiv.org/abs/1610.09027)
 220 09027.
- 221 [9] A. Santoro, D. Raposo, D. G. T. Barrett, M. Malinowski, R. Pascanu, P. Battaglia, and T. Lil-
 222 licrap. A simple neural network module for relational reasoning. pages 1–16, 2017. ISSN
 223 21607516. doi: 10.1109/WACV.2017.108. URL <http://arxiv.org/abs/1706.01427>.
- 224 [10] J. Weston, A. Bordes, S. Chopra, A. M. Rush, B. van Merriënboer, A. Joulin, and T. Mikolov.
 225 Towards AI-Complete Question Answering: A Set of Prerequisite Toy Tasks. 2015. ISSN
 226 03787753. doi: 10.1016/j.jpowsour.2014.09.131. URL [http://arxiv.org/abs/1502.](http://arxiv.org/abs/1502.05698)
 227 05698.
- 228 [11] J. Weston, S. Chopra, and A. Bordes. Memory Networks. pages 1–15, 2015. ISSN 1098-7576.
 229 doi: v0. URL <https://arxiv.org/pdf/1410.3916.pdf>.