

# Machine Comprehension Task with BiDirectional Attention Flow on SQuAD Dataset

**Shim-Young Lee**  
Department of Statistics  
Stanford University  
*sylee1@stanford.edu*

## Abstract

In this default project (assignment 4) for the course CS224n Natural Language Processing with Deep Learning, a challenging NLP task called machine comprehension was explored in the form of question answering. The SQuAD dataset provides the contexts and corresponding questions, and the task is to correctly highlight, or mark the beginning and the end of, a word sequence in the context paragraph that answers the question. I augmented the baseline model with the attention mechanism of the BiDirectional Attention Flow and Bidirectional LSTMs in the contextual embed layer and modeling layer of the neural network. A single model achieved THIS F1 and THIS EM.

## 1 Introduction

The goal of this project is to explore a natural language processing (NLP) task of machine comprehension in the form of question answering, where the level of comprehensive capability of a machine is assessed by evaluating the machine generated answers given a context paragraph and a question about the context. Machine comprehension is a challenging task, as the machine must be able to extract information from a relatively long sequence of words and represent, or summarize that information in a form that will allow the machine to use when answering a question. It shares many challenges of machine translation as it must understand the relationship between words and sentences; however, instead of mapping a word sequence in one language to another, it also must extract useful information in both the context and the question and model the complex interaction between them in order to come up with a correct answer. The “answer” in this specific task with SQuAD dataset[4] is given by marking the start and the end position of a word sequence in the context paragraph.

In this project, a baseline model with GRU[1] cells in the encoding layer and a basic attention mechanism was provided. The task was to replace or augment this provided baseline model in a single or multiple modules of the network architecture to achieve better performance. I chose to implement Bidirectional Attention Flow (BiDAF) [5] whose architecture has the similar modularization/hierarchy as the baseline model. BiDAF is one of the highest performing model for this task one SQuAD dataset.

41 **2 Related Work**

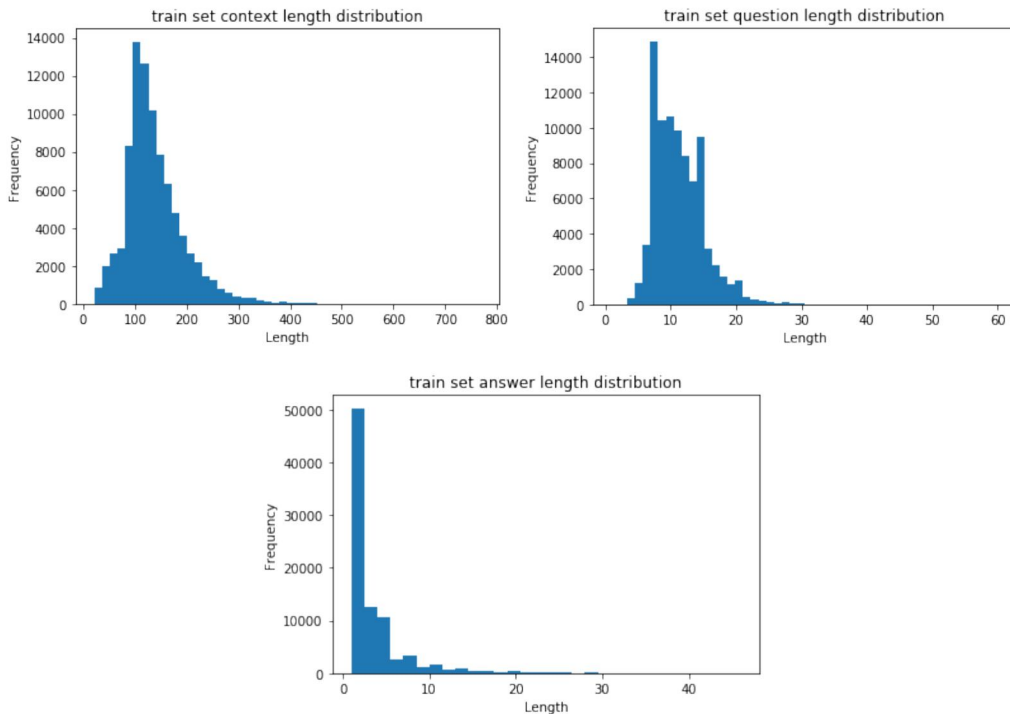
42 As previously stated, the related work that was most influential to this project is the  
43 Bidirectional Attention Flow model, whose attention mechanism was implemented here. The  
44 general organization of components or layers other than the attention layer is also very similar  
45 for the BiDAF and the version I implemented for this project. The full BiDAF model has other  
46 extensions such as character-level CNN embedding layer.

47 Other important model that performed well on SQuAD dataset is Dynamic Coattention  
48 Network by Xiong *et al*[7]. Similarly to the BiDAF model, one of the major contributions of  
49 the Coattention Network is its attention mechanism in the Coattention layer, which uses a two-  
50 way attention between the context and the question, as well as the second-level attention  
51 computation that attends over attention outputs.

52  
53  
54 **3 Dataset and Features**

55 The dataset used for the question answering task here is the Stanford Question Answering  
56 Dataset (SQuAD). SQuAD is a reading comprehension dataset; it consists of about 100,000  
57 context paragraph-question pairs extracted from Wikipedia articles as well as the answers to  
58 the questions, built through crowdsourcing. From the perspective of a model, each input is a  
59 pair of context paragraph and a question about that paragraph, and the goal is to answer the  
60 question correctly. The training set has roughly 85,000 context-question pairs, and the dev set  
61 has about 10,000 pairs. The test set is similar to the dev set in terms of size, and I did not have  
62 direct access to the test set, and it was used only for the final evaluation/submission.

63 To make the assessment of the correctness of the machine generated answers easy, the answers  
64 are always taken directly from the context paragraph. Specifically, the answer is specified by  
65 marking the start and the end positions of the word sequence in the context paragraph, so the  
66 question answering system needs not generate a new word sequence. The accuracy of the  
67 answers are evaluated simply by comparing these machine generated start and end positions  
68 with the answers provided by the dataset.



69

70

71 Figure 1: Histogram of sequence lengths for context paragraphs, questions and answers in  
72 training set.

73 Table 1: 95 and 99 percentile of the sequence lengths for context paragraphs, questions and  
74 answers in training set.

75

<b>Word sequence type</b>	<b>95 percentile length</b>	<b>99 percentile length</b>
Context	245	325
Question	18	23
Answer	10	18

76

77 The training and inference with the model involve tensor multiplications, so the maximum  
78 sequence lengths of a context paragraph and a question are hyperparameters to be set. To  
79 determine the cut-off length for each effectively, histograms of sequence lengths were used.  
80 The figure 1 and table 1 summarize the distributions of lengths for context paragraphs,  
81 questions and answers. Based on these statistics, the maximum length for context paragraphs  
82 was set to 300 and for questions 25.

83

#### 84 **4 Model Architecture**

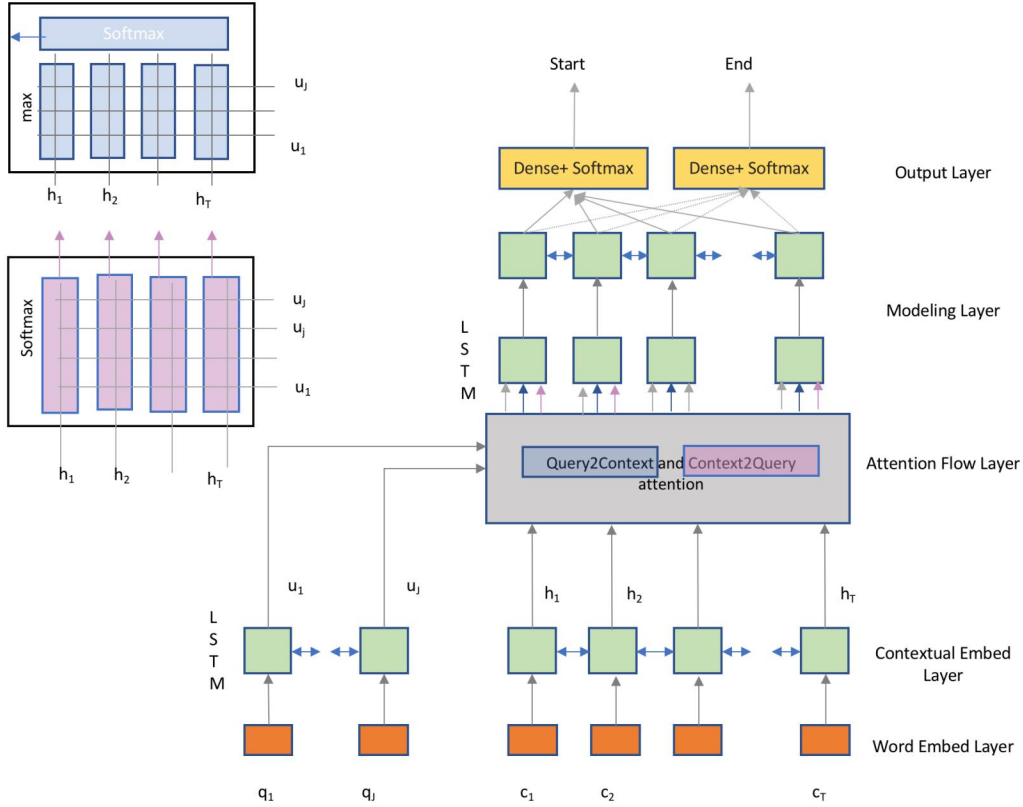
85 First, I provide the higher level description of the modularization and hierarchy of the model  
86 architectures that all models tested follow.

87 Both the baseline model and the BiDAF models consist of three major components, or layers.  
88 The first layer (besides the very first word embedding layer) is the RNN encoder layer, which  
89 takes in the word level embedding vectors of the sentence and uses some RNN encoder (GRU  
90 or LSTM) to encode the word sequence into hidden states. In this layer, context paragraph and  
91 question are encoded with the same RNN cell, so the weights are shared.

92 The next layer is the attention layer. The interaction/similarity between context and question  
93 hidden states are computed, with dot-product in the simplest case, and the score is turned into  
94 a distribution by applying softmax to the scores. The distribution is then used as a weight in  
95 the weighted sum of either context or question hidden states, which is blended together to be  
96 fed into following layers.

97 The layer following the attention layer is the output layer, where the blended representation  
98 of reweighted context and question vectors are fed into the usual fully connected layers.  
99 Specifically, affine transformation followed by nonlinear activation (ReLU was used for all  
100 models discussed) is applied to the blended representation, and the output once again goes  
101 through another affine transformation and then the softmax is applied to turn these logits into  
102 a probability distribution. Note that the probability distribution here is the probability of each  
103 word in the context paragraph being the start/end position of the highlight span, so two output  
104 layers are needed, each for computing the distribution for the start position, and the end  
105 position. The final prediction is made simply by taking the argmax of the probability  
106 distribution.

107 Next, the implementation details for each component of respective models are discussed.



108  
109

110 Figure 2: Model Architecture diagram for BiDirectional Attention Flow model

111  
112  
113

#### 4.1 Baseline model

114 **Embedding Layer:** The GloVe [3] word embedding vectors of size 100 was used.  
 115 **RNN Encoder Layer:** A bidirectional GRU was used for both context and the question. An  
 116 embedded sequence (GloVe) of context and query vectors  $\mathbf{c} = [c_1, c_2, \dots, c_T]$  and  
 117  $\mathbf{q} = [q_1, q_2, \dots, q_J]$  are fed into the RNN cell and the representation matrix  $\mathbf{H} \in \mathbb{R}^{T \times 2h}$  and  
 118  $\mathbf{U} \in \mathbb{R}^{J \times 2h}$  each for context and question are output ( $2h$  since bidirectional, where  $h$  is the  
 119 size of a single hidden state).  
 120 **Attention Layer:** A simple dot-product attention was used: similarity matrix  $\mathbf{S} = \mathbf{H}_t \mathbf{U}_j^T$ ,  
 121 attention weight vector  $\mathbf{a}_t = \text{softmax}(\mathbf{S}_{t,:}) \in \mathbb{R}^J$  was computed and used to reweight the  
 122 question vectors, yielding  $\tilde{\mathbf{U}}_t = \sum_j \mathbf{a}_{tj} \mathbf{U}_j^T \in \mathbb{R}^{T \times 2h}$ .  
 123 **Output Layer:** A simple output layer described above was used.

124  
125  
126

#### 4.2 Bidirectional Attention Flow

127 **Embedding Layer:** The GloVe word embedding vectors of size 100 was used. GloVe vectors  
 128 of higher dimension were not used nor tested due to the training speed and the GPU memory  
 129 limitation.

130 **RNN Encoder Layer:** Two layers of bidirectional LSTM was used. This is the contextual  
 131 embed layer in the figure 2 (as the authors of BiDAF called it).

132

133 **Attention Layer:** The similarity matrix  $\mathbf{S}_{ij} = \mathbf{w}_{sim}^T [\mathbf{H}_{t_i}, \mathbf{U}_{j_i}; \mathbf{H}_{t_i} \circ \mathbf{U}_{j_i}] \in \mathbb{R}$  was computed, then  
 134 the Context2Question attention,  $\alpha^t = \text{softmax}(\mathbf{S}_{t,:}) \in \mathbb{R}^J$ ,  $\tilde{\mathbf{U}}_t = \sum_{j=1}^J \alpha_j^t \mathbf{U}_{j_i}$  and the  
 135 Question2Context attention  $\mathbf{m}_t = \max_j \mathbf{S}_{t,j} \in \mathbb{R}$ ,  $\beta = \text{softmax}(\mathbf{m}) \in \mathbb{R}^N$ ,  $\mathbf{h} = \sum_{t=1}^T \beta_t \mathbf{H}_{t_i}$ ,  
 136 and  $\tilde{\mathbf{H}}$  is constructed by tiling  $\mathbf{h}$  to match the dimension of the final blended representation,  
 137  $\mathbf{B} = [\mathbf{H}; \mathbf{U}; \mathbf{H} \circ \tilde{\mathbf{U}}, \mathbf{H} \circ \tilde{\mathbf{H}}]$

138 **Modeling Layer:** This is a layer not described in the higher lever overview in the beginning  
 139 of the section. Instead of feeding the blended representation right into the output layer, two  
 140 layers of bidirectional LSTM were used to encode the representation.

141 **Output Layer:** A simple output layer described above (same as the baseline model) was used.  
 142 The figure 2 summarizes the overall architecture and flow of computation.

143

144

## 145 4 Result and Discussion

146

147 In assessment of the performance of question answering task, F1 score and Exact Match  
 148 (EM) score are used. F1 score is the harmonic mean of precision and recall, and Exact Match  
 149 is a simple accuracy score that computes the proportion of predictions that exactly match the  
 150 answer span (given by the start and the end positions).

151

152 The baseline model achieves F1 score of 43.027 and the EM score of 34.134 in the dev set  
 153 with about 16k iterations of training step with batch size 100.

154

155 The first augmented model trained was the BiDAF model without the modeling layer  
 156 (BiDAF1), that is, the baseline model with replaced attention mechanism (one layer of GRU  
 157 RNN Encoder layer, BiDAF attention layer, and output layer). This model achieved the F1  
 158 score of 48.77 and EM score of 35.47 in the dev set with over 20k iterations of training steps  
 159 with batch size 100. The improvement in performance was fairly minor given the increase in  
 160 the complexity of the attention layer.

161

162 Finally, the BiDAF model with one layer of bidirectional LSTM and two layers of  
 163 bidirectional LSTM in the contextual embed layer and the modeling layer respectively, and  
 164 the BiDAF attention mechanism (BiDAF2), achieved the F1 score of 73.115 and EM score  
 165 of 63.084 in the dev set, and F1 score of 73.553 and EM score of 64.219 in the test set with  
 166 only about 12k iterations of batch size 64.

167

168

169

170 Table 2: The F1 and EM scores of each model.

171

172

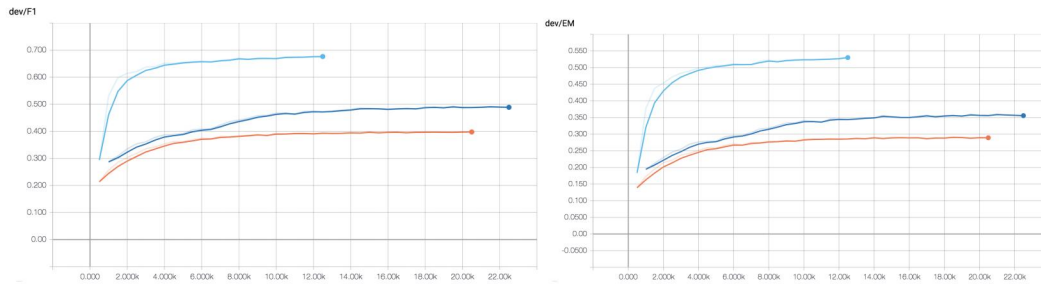
Model	F1	EM
Baseline(dev)	43.027	34.134
BiDAF1(dev)	48.77	35.47
BiDAF2(test)	73.553	64.219

173

174

175

176



177

178

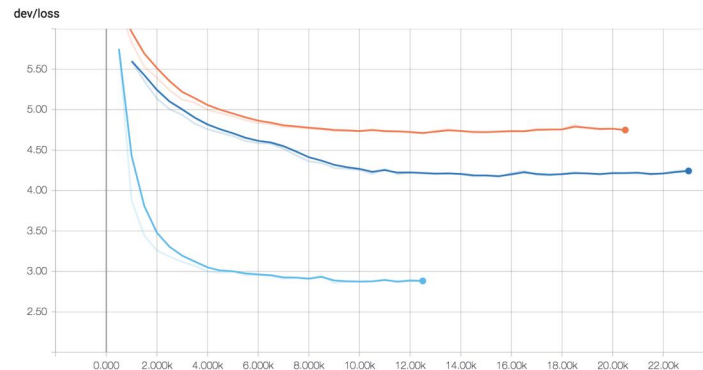
179

Figure 3: Dev set F1 score (left) and Dev set EM score (right) as a function of iterations.

Orange: baseline, navy: GRU BiDAF model with no modeling layer,

180

sky blue: LSTM BiDAF with a modeling layer



181

182

Figure 4: Dev set loss as a function of iterations.

183

Orange: baseline, navy: GRU BiDAF model with no modeling layer,

184

sky blue: LSTM BiDAF with a modeling layer

185

186

187 The figures above show the dev set evaluation of F1 and EM scores as well as the loss of  
 188 three models as a function of iterations. (I did not save the loss and F1/EM scores separately,  
 189 so I took the figures directly from the TensorBoard visualization tool.)

190 All three plots show that BiDAF2 achieves loss and performance in only a few thousand  
 191 iterations the other two models couldn't achieve in more than twice the training steps. Also  
 192 keep in mind that the batch size was smaller for BiDAF2, due to the GPU memory constraint.

193 The overwhelming difference in the performance of BiDAF1 and BiDAF2 suggests that  
 194 BiDAF's attention mechanism alone cannot improve the performance of the baseline model  
 195 without the modeling layer. Even with the increase in model complexity in BiDAF2 due to the  
 196 change in the RNN encoder from GRU to LSTM and two layers of modeling layer, overfitting  
 197 was much more severe in BiDAF1, where the training F1 and EM scores went up as high as  
 198 over 85. On the other hand, for BiDAF1, the F1 and EM scores for training set and dev set  
 199 remained roughly similar throughout the training, and the model performed better on the test  
 200 set than on the dev set. This might be due to the decrease in the size of hidden states for all  
 201 RNN encoders (was forced to do so due to the memory constraint).

202

## 203 5 Conclusion and Future Work

204 For this default project, I implemented some of the components of the BiDirectional Attention  
 205 Flow to augment the baseline model, and was able to achieve improvement of about 25 in F1  
 206 score and 30(%) in EM score. The more complex attention mechanism, however, improves the  
 207 model significantly only when the contextual embed layer and the modeling layer are  
 208 sufficiently complex.



209 One possible, obvious extension to the current model is the character-level CNN, which uses  
210 the character embedding to represent the characters and provides more subtle and refined  
211 interaction of word sequences in character level in addition to the word level embedding.

## 212 **References**

213

214 [1] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio,  
215 “Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation,” ArXiv  
216 eprints, Jun. 2014.

217 [2] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Comput.*, vol. 9, no. 8,  
218 pp.1735-1780, Nov. 1997.[Online]. Available: <http://dx.doi.org/10.1162/neco.1997.9.8.1735>

219

220 [3] Pennington, J., Socher, R., & Manning, C. D. “GloVe: Global Vectors for Word Representation,” 1532–  
221 1543. Retrieved from <http://www.aclweb.org/anthology/D14-1162>, 2015

222

223 [4] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, “SQuAD: 100,000+ Questions for Machine  
224 Comprehension of Text,” ArXiv e-prints, Jun. 2016.

225

226 [5]M. Seo, A. Kembhavi, A. Farhadi, and H. Hajishirzi, “Bidirectional attention flow for machine  
227 comprehension,” arXiv preprint arXiv:1611.01603, 2016.

228

229 [6] Sutskever, I., Vinyals, O., and Le, Q.. “Sequence to sequence learning with neural networks,” In *Advances*  
230 *in Neural Information Processing Systems*. 2014

231

232 [7] C. Xiong, V. Zhong, and R. Socher, “Dynamic Coattention Networks For Question Answering,” ArXiv e-  
233 eprints, Nov. 2016.