

# SQuAD Model Exploration: BiDAF and Input Features

**Ben Barnett and Jeffrey Chen**  
Department of Computer Science  
Stanford University  
{ben.barnett, jchen623}@stanford.edu

## Abstract

In this paper, we explore various deep learning techniques and approaches to implementing a reading comprehension model using the Stanford Question Answering Dataset (SQuAD), a reading comprehension dataset that includes a set of Wikipedia articles, and crowdsourced pairs of 100k+ questions and answers to those corresponding articles. Specifically, each corresponding answer is taken directly from each question, meaning each answer is an exact excerpt or “span” of the original context. Through a model inspired by “Bidirectional attention flow for machine comprehension” [1] and “Reading Wikipedia to answer open-domain questions” [2], we found that adding BiDAF to an RNN model with additional “word embedding” and “exact match” input features improved upon our original baseline model, with an F1 score of .43 and an EM score of .30.

## 1 Introduction

### 1.1 Objective

Our task, is to use various deep learning techniques to develop a model that performs well on the Stanford Question Answering Dataset (SQuAD). Machine reading comprehension and question answering is a complex, yet fundamental problem in Natural Language Processing with endless possible approaches. Our goal is to be able to take a context paragraph and question, and return an answer. For example, consider the context and question below:

Context: “during his time at his lab , tesla observed unusual signals from his receiver which he concluded may be communications from another planet . he mentioned them in a letter to reporter julian hawthorne at the philadelphia north american on 8 december 1899 and in a december 1900 letter about possible discoveries in the new century to the red cross society where he referred to messages " from another world " that read " 1 ... 2 ... 3 ... " . reporters treated it as a sensational story and jumped to the conclusion tesla was hearing signals from mars . he expanded on the signals he heard in a 9 february 1901 collier's weekly article " talking with planets " where he said it had not been immediately apparent to him that he was hearing " intelligently controlled signals " and that the signals could come from mars , venus , or other planets . it has been hypothesized that he may have intercepted marconi 's european experiments in july 1899—marconi may have transmitted the letter s ... in a naval demonstration , the same three impulses that tesla hinted at hearing in colorado—or signals from another experimenter in wireless transmission.”

42 **Question:** “To what did tesla attribute the unknown signals his radio received ?”  
43 Given this context and question, our model should return the answer “communications from  
44 another planet.” Notice that this answer is a “span” or direct excerpt from the given passage.

45  
46 **1.2 Performance Metrics**

47 The passage above is a direct data point from SQuAD, and the performance on the dataset is  
48 measured with 2 metrics: F1 and Exact Match (EM) scores. F1 is defined as  
49  $2 \times \text{precision} \times \text{recall} / (\text{precision} + \text{recall})$ . EM, or “exact match” on the other hand is a binary  
50 measurement of whether the predicted value matches the true answer exactly. For reference,  
51 human performance on this dataset is an F1 score of .86 [3]. Without sophisticated  
52 knowledge of linguistics and semantics, we attempt to use deep neural networks to perform  
53 as well as possible on SQuAD.

54  
55 **1.2 Related Works**

56 There has been an extensive amount of research regarding SQuAD and an even greater  
57 concentration on Machine Reading Comprehension as a whole. However there are a few  
58 notable papers in which we gathered inspiration for our own improvements. “Bidirectional  
59 attention flow for machine comprehension” [1] describes the appropriate architecture for  
60 implementing BiDAF. This paper also describes many other improvements to be made to the  
61 baseline model like Logistic Regression, Co-attention, and Fine-Grained Gating, but we  
62 believed that BiDAF was the primary modification that resulted in their model’s high  
63 performance. We also added additional input features with the inspiration of “Reading  
64 Wikipedia to answer open-domain questions” [2]. This paper describes significant  
65 improvements from adding a few other additional features to the word embeddings. They  
66 found that adding these features below specifically improved F1 scores most significantly.

- 67 1. Representing if each context token appeared inside the question
- 68 2. Adding Part-of-Speech tags and Named Entity types to each context
- 69 3. Using the word embedding for each context to attend to the word embeddings for  
70 the question

71  
72 **2 Approach**

73 The general architecture of the final model is the combination of several modifications to the  
74 baseline model provided by the CS224N course staff. The following sections outline the  
75 architecture of the baseline model and the individual modifications to the baseline model  
76 explored in this project.

77  
78 **2.1 Baseline Model**

79 The baseline model consists of three components: a RNN encoder layer that encodes  
80 contexts and questions into hidden states, an attention layer that combines context and  
81 question hidden states into a single representation, and an output layer, which outputs the  
82 starting and end location of the predicted answer span within the context.

83  
84 **2.1.1 RNN Encoder Layer**

85 Every context and question word is mapped to its corresponding d-dimensional, pre-trained  
86 GLoVe embedding. These fixed embeddings are then fed into a 1-layer bidirectional GRU,  
87 which produces a sequence of h-dimensional forward hidden states and h-dimensional  
88 backward hidden states for both contexts and questions:

$$\begin{aligned} \{\vec{c}_1, \vec{c}_1, \dots, \vec{c}_N, \vec{c}_N\} &= \text{biGRU}(\{\mathbf{x}_1, \dots, \mathbf{x}_N\}) \\ \{\vec{q}_1, \vec{q}_1, \dots, \vec{q}_M, \vec{q}_M\} &= \text{biGRU}(\{\mathbf{y}_1, \dots, \mathbf{y}_M\}) \end{aligned}$$

89  
90 Figure 1: Output of Bidirectional GRU within RNN Encoding Layer

91

92 Then, the forward and backward hidden states are concatenated, which comprises the output  
93 hidden states for each corresponding question and context word from the encoding layer:

$$\begin{aligned} \mathbf{c}_i &= [\vec{\mathbf{c}}_i; \overleftarrow{\mathbf{c}}_i] \in \mathbb{R}^{2h} \quad \forall i \in \{1, \dots, N\} \\ \mathbf{q}_j &= [\vec{\mathbf{q}}_j; \overleftarrow{\mathbf{q}}_j] \in \mathbb{R}^{2h} \quad \forall j \in \{1, \dots, M\} \end{aligned}$$

94

95 Figure 2: Final Hidden Representations Output from RNN Encoder Layer

96

### 97 2.1.2 Attention Layer

98 The baseline attention layer consists of a basic dot-product attention in which context hidden  
99 states attend to the question hidden states. This layer outputs an  $N$  long sequence of blended  
100 representations wherein each representation corresponds to one context word and  $N$  is the  
101 context length.

102

### 103 2.1.3 Output Layer

104 Each blended representation,  $b_i$ , is fed through a fully connected layer followed by a ReLU  
105 activation function, which is then passed through a linear layer that computes the starting  
106 score of each context word:

$$\begin{aligned} \mathbf{b}'_i &= \text{ReLU}(\mathbf{W}_{FC}\mathbf{b}_i + \mathbf{v}_{FC}) \in \mathbb{R}^h \quad \forall i \in \{1, \dots, N\} \\ \text{logits}_i^{\text{start}} &= \mathbf{w}_{\text{start}}^T \mathbf{b}'_i + u_{\text{start}} \in \mathbb{R} \quad \forall i \in \{1, \dots, N\} \end{aligned}$$

107

108

109 Figure 3: Computation of Context Blended Representations

110

111 To obtain the probability distribution of start words,  $\text{logits}_i^{\text{start}}$  is passed through a final  
112 softmax layer. Computing the predicted end word is a parallel process, with the only  
113 difference being the weight and bias terms  $w_{\text{end}}$  and  $u_{\text{end}}$  replacing  $w_{\text{start}}$  and  $u_{\text{start}}$ .

114

### 115 2.1.4 Loss and Predictions

116 The loss function for a single example is computed as the sum of the cross-entropy loss for  
117 start and end locations, using  $i_{\text{start}}$  and  $i_{\text{end}}$  for the true starting and ending word locations:

$$\text{loss} = -\log p^{\text{start}}(i_{\text{start}}) - \log p^{\text{end}}(i_{\text{end}})$$

118

119 Figure 4: Loss Calculation for a Single Context

120

121 During training, loss is minimized across the average of each batch, using an Adam  
122 optimizer.

123 Note that during testing, rather than making start and end word predictions based on their  
124 respective probability distributions, the start and end words with the highest respective  
125 scores are output as predictions.

126

## 127 2.2 Modification 1: Bidirectional Attention Flow

128 Bidirectional Attention Flow (“BiDAF”) is a modification to the baseline model that  
129 improves the attention layer by allowing attention to flow in both directions. In the baseline  
130 model, the context attends to the question. However, in BiDAF, the question attends to the  
131 context as well, resulting in a higher dimensional blended representation of each context  
132 word whose features are hopefully more informative for the output layer of the model. The  
133 C2Q component is computed similarly to the baseline model.

134 All equations and formulas for computation can be found on the final project handout,  
 135 though the key part of BiDAF is the altered blended representation of each context word,  
 136 which contains twice as many features at the end of the layer (for a total of  $8h$ -  
 137 dimensionality). The blended representation for a single word is as follows:

$$b_i = [c_i; a_i; c_i \circ a_i; c_i \circ c'] \in \mathbb{R}^{8h} \quad \forall i \in \{1, \dots, N\}$$

138

139 Figure 5: Final Blended Representation from BiDAF Attention Layer

140

141 wherein  $c_i$  is the context word's embedding and  $a_i$  is the C2Q attention output—a metric of the  
 142 context word to all question words.  $c'$  is the Q2C attention output that is new to the model from  
 143 the previous baseline implementation which used only C2Q, and is computed as follows:

144

$$S_{ij} = w_{\text{sim}}^T [c_i; q_j; c_i \circ q_j] \in \mathbb{R}$$

145

$$m_i = \max_j S_{ij} \in \mathbb{R} \quad \forall i \in \{1, \dots, N\}$$

$$\beta = \text{softmax}(m) \in \mathbb{R}^N$$

$$c' = \sum_{i=1}^N \beta_i c_i \in \mathbb{R}^{2h}$$

146

147 Figure 6: Calculation of  $c'$ , the Q2C Attention Output during BiDAF

148

149 In the equations above,  $S_{ij}$  is the similarity score between  $c_i$  and  $q_j$ . In order to compute similarities  
 150 efficiently, a single similarity matrix  $S \in \mathbb{R}^{N \times M}$  is computed once for every context, question pair  
 151 in the batch, whose entries represent the similarity between each combination of word pairs within  
 152 the context and question.

153

### 154 2.3 Modification 2: Additional Input Features

155 In our final model, we concatenated additional input features to the output produced in the  
 156 baseline encoding layer. Two features were appended to each context word embedding:

157 1. A positive integer metric, representing whether the context word appears in the  
 158 question.

159 2. A scalar representing the minimum Euclidean distance from the context word's  
 160 GLoVe embedding and all question word embeddings.

161 More formally, given a context word's GLoVe embedding,  $c\text{-emb}_i \in \mathbb{R}^H$  (the GLoVe  
 162 embedding hidden size), the context word's final embedding output from the model's  
 163 encoding layer *before passing through the bidirectional GRU* is as follows:

$$164 [c\text{-emb}_i; a; b] \in \mathbb{R}^{H+2}$$

165 where  $a$  is the number of times that the context word appears in the question (most often just  
 166 a 0 or 1) and

$$167 b = \min_j \| c\text{-emb}_i - q\text{-emb}_j \| \quad j \in \{1, \dots, M\}.$$

168 The intuition behind these features is that often when trying to find the answer to a question  
 169 in a passage, words in the passage that appear in the question are likely to be near or part of  
 170 the answer. This is also part of why basic dot-product attention is so effective, as similarities  
 171 between context and question words are often a strong indicator of where the answer lies.

172 Incorporating this technique/intuition into the encoding layer seemed like a good way to  
173 directly ‘encode’ the importance of similarity between a context word and its existence  
174 within the corresponding question.

175

### 176 **3 Experiments**

177

#### 178 **3.1 Dataset**

179 SQuAD is a crowdsourced reading comprehension dataset containing context, question,  
180 answer triplets from Wikipedia entries [3]. The dataset used includes about 100,000  
181 questions in total, though note that there are fewer distinct contexts, as multiple questions  
182 and answers may come from the same context. More information about SQuAD can be found  
183 at <https://rajpurkar.github.io/SQuAD-explorer/explore/1.1/dev/>.

184

#### 185 **3.2 Evaluation Metrics**

186 The primary performance metrics used were F1 and EM scores, whose definitions can be found in  
187 the introduction section.

188

#### 189 **3.3 Results**

190

Table 1: Baseline

191

Model configuration	Default baseline model as described above
Learning rate	0.001
Batch size	32
Num iterations reached	40k
Training time	15 hrs.

192

193

Table 2: Baseline + BiDAF

194

Model configuration	Default baseline model as described above with the exception of the model’s attention layer, which was substituted with a BiDAF implementation
Learning rate	0.001
Batch size	32
Num iterations reached	12k
Training time	28 hrs.

195

196

Table 3: Baseline + BiDAF + additional input features

197

Model configuration	Default baseline model with modified RNN encoder layer to include additional input features as well as substitution of the default attention layer for a BiDAF attention layer
Learning rate	0.001
Batch size	32
Num iterations reached	15k
Training time	40 hrs.

198

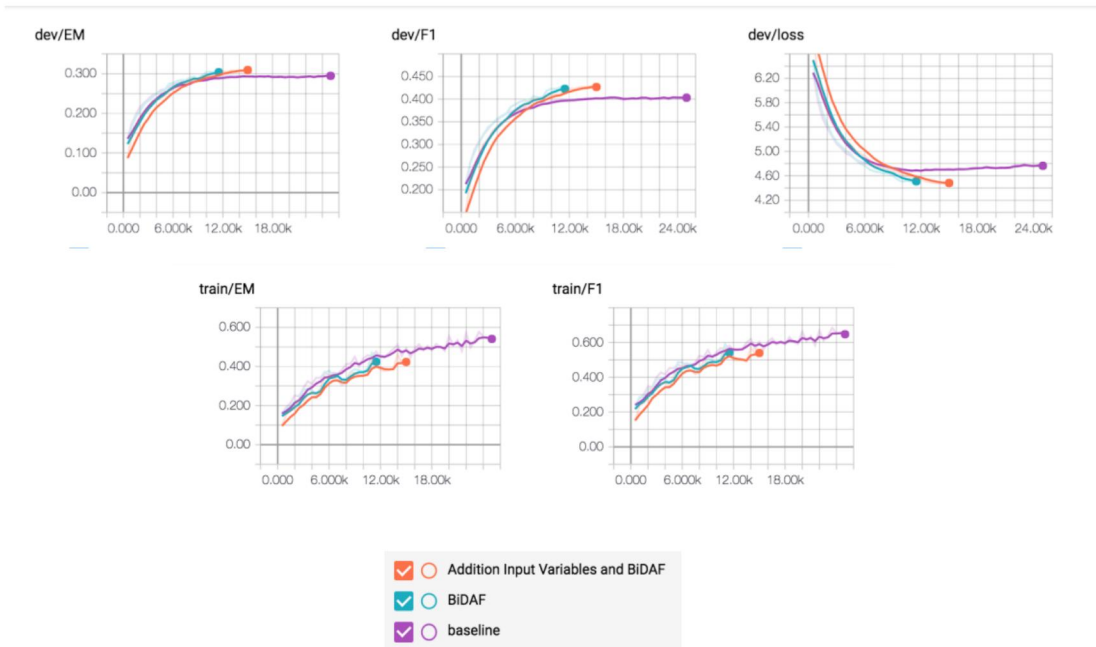


Figure 7: Various model results overlaid over iterations run

205 From Figure 7 above, we can see that BiDAF with our additional input features achieves the  
 206 highest F1 score, plateauing at .43. This model also performs the best in terms of with a  
 207 score of .31. By analyzing individual example answer spans that our model predicted, we  
 208 were able to identify specific strengths and weaknesses of our implementation through the  
 209 following examples.

#### EXAMPLE 1

211 CONTEXT: a piece of paper was later found on which Luther had written his last statement.  
 212 The statement was in Latin , apart from " we are beggars , " which was in German.

214 QUESTION: what was later discovered written by luther ?

215 TRUE ANSWER: his last statement

216 PREDICTED ANSWER:

217 F1 SCORE ANSWER: 0.000

218 EM SCORE: False

#### EXAMPLE 2

220 CONTEXT: after leaving edison 's company tesla partnered with two businessmen in 1886 ,  
 221 robert lane and benjamin vail , who agreed to finance an electric lighting company in tesla 's  
 222 name , tesla electric light & manufacturing . the company installed electrical arc light based  
 223 illumination systems designed by tesla and also had designs for dynamo electric machine  
 224 commutators , the first patents issued to tesla in the us .

225 QUESTION: what was produced at tesla 's company ?

226 TRUE ANSWER: dynamo electric machine commutators

227 PREDICTED ANSWER: electrical arc light based illumination systems

228 F1 SCORE ANSWER: 0.000

229 EM SCORE: False

230

### EXAMPLE 3

231 CONTEXT: the ipcc receives funding through the ipcc trust fund , established in 1989 by the  
232 united nations environment programme ( unep ) and the world meteorological organization (   
233 wmo ) , costs of the secretary and of housing the secretariat are provided by the wmo , while  
234 unep meets the cost of the depute secretary . annual cash contributions to the trust fund are  
235 made by the wmo , by unep , and by ipcc members ; the scale of payments is determined by  
236 the ipcc panel , which is also responsible for considering and adopting by consensus the  
237 annual budget . the organisation is required to comply with the financial regulations and  
238 rules of the wmo .

239 QUESTION: who funds the ipcc 's deputy secretary ?

240 TRUE ANSWER: united nations environment programme

241 PREDICTED ANSWER: united nations environment programme ( unep ) and the world  
242 meteorological organization ( wmo ) , costs of the secretary and of housing the secretariat  
243 are provided by the wmo

244 F1 SCORE ANSWER: 0.320

245 EM SCORE: False

246

### EXAMPLE 4

247 CONTEXT: in june 1978 , arledge created the newsmagazine 20/20 ; after its first episode  
248 received harshly negative reviews , the program – which debuted as a summer series , before  
249 becoming a year-round program in 1979 – was immediately revamped to feature a mix of in-  
250 depth stories and interviews , with hugh downs appointed as its anchor ( later paired  
251 alongside his former today colleague barbara walters ) . in february 1979 , abc sold its  
252 recording division to mca inc. for \$ 20 million ; the label was discontinued by march 5 of  
253 that year , and all of its 300 employees were laid off ( the rights to the works of abc records  
254 and all of mca 's other labels have since been acquired by universal music group ) .

255 QUESTION: when was the newsmagazine 20/20 first created ?

256 TRUE ANSWER: june 1978

257 PREDICTED ANSWER: june 1978

258 F1 SCORE ANSWER: 1.000

259 EM SCORE: True

260

261 From example 1, you can see that our model fails to make any prediction whatsoever. This is  
262 because the span start-index was predicted to be after the end-index. Furthermore, in  
263 example 3, we can see that the correct answer is only 4 words long, whereas our model  
264 predicted the answer to have the gold truth starting index, but to span 25 words. Although  
265 this answer is technically correct (in that it contains the correct answer), it received a low F1  
266 score because it was much longer than necessary. From these examples, among others, we  
267 determined that regulating span would be an important modification to make in the future.  
268 Specifically, ensuring that span start indices were never predicted to be after end indices and  
269 that span lengths are always under a predetermined size would not only improve accuracy  
270 and F1 scores, but also the efficiency of our model.

271 Example 2 highlights another shortcoming of our model. The question asks “what was  
272 produced at tesla 's company?” The true answer is “dynamo electric machine commutators”,  
273 but our model instead predicts the answer to be “electrical arc light based illumination  
274 systems.” We attribute this error to the additional input feature which adds a similarity  
275 metric between GloVe contexts and questions to each word embedding. This input feature  
276 directly affects our model in this example because the question specifically asks “what was  
277 produced...?” and our model likely uses the fact that “designed” has a closer GloVe  
278 embedding to “produced” than “installed.”

279 However, this same feature seems to help our model perform better on questions inquiring  
280 about numbers and dates. In example 4, our model scores a perfect F1 of 1.00. This is likely  
281 because numbers exist in a very specific space within the entire GloVe embedding space,  
282 which results in greater certainty in our model when searching for numbers and context  
283 tokens similar to numbers. This ultimately improves the model’s ability to find the correct  
284 answer span when prompted for an answer to contain dates or numbers.

285

## 286 **5 Conclusion**

287 Bidirectional attention flow and additional input features both resulted in improvements  
288 from the baseline model. However, in analyzing the results of example sentences (and the  
289 model’s predictions), we noticed that there still seems to be a large opportunity for “easy”  
290 model improvement by training the model to more heavily consider the answer span in the  
291 output layer of the model before simply predicting start and end words independently. For  
292 example, the final model often outputs no answer at all because its predicted end word is  
293 prior to the predicted start word; preventing this alone by requiring the start word to be  
294 before the end word may alone cause a large model improvement without any significant  
295 increase in the computational expense of training. Further, it’s possible that improvements  
296 would be visible if the model were to consider joint probabilities between the start and end  
297 word as opposed to calculating the most likely start and end words independently.

298 One other broad takeaway outside of the primary objective of investigating the best possible  
299 SQuAD model is the importance of the model’s training efficiency; although our final model  
300 outperformed the baseline model in all evaluation metrics, it took roughly 7 times slower to  
301 train, making it expensive to tune hyperparameters. In the future, given a perhaps even larger  
302 dataset, it would be vital to find ways to improve the efficiency of the model, especially the  
303 implementation of BiDAF.

304

## 305 **6 References**

306 [1] Minjoon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. Bidirectional  
307 attention flow for machine comprehension. arXiv preprint arXiv:1611.01603, 2016

308

309 [2] Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading wikipedia to  
310 answer open-domain questions. arXiv preprint arXiv:1704.00051, 2017.

311

312 [3] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100, 000+  
313 questions for machine comprehension of text. CoRR, abs/1606.05250, 2016.