

# Deep Learning for Question Answering on the SQuAD

1                   **Ramtin Keramati**      **Kian Katanforoosh**      **Reza Takapoui**  
2                   PhD ICME                   MS MS&E                   PhD EE  
3                   keramati@stanford.edu   kian@stanford.edu   takapoui@stanford.edu  
4

## Abstract

5                   This paper describes our project for the Spring 2017 session of CS224N.  
6                   Question answering has been a hard problem for a long time, and it was  
7                   originally hard to apply cutting-edge Deep Learning algorithm to this  
8                   problem because of lack of data. In 2016, the Stanford Question Answering  
9                   Dataset (SQuAD) was published. It boosted the application of deep learning  
10                  algorithm to Question Answering problems.

11  
12                  In fact, we will see in this paper two different algorithms to tackle this  
13                  problem. We will start with a Baseline algorithm that relies on  
14                  Bidirectional-LSTMs, and later boost its performances by adding an  
15                  attention mechanism.

16

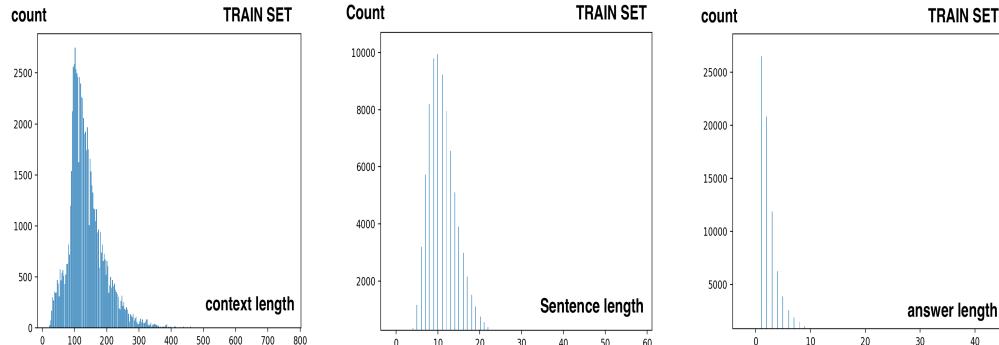
## 1      Data Analysis

17                  Before jumping to the architecture of our model, we would like to see what our dataset look  
18                  like. It will help us take the right decision and make hypothesis on our model's expected  
19                  behavior.

20

### 22      2.1 Sequences length

23                  First, let us derive some interesting plots on the frequency distribution of sequences  
24                  (questions = sentence, context=paragraph and answers) length on the training set.



25

26                  By looking at these plots, we can assert that:

- 27                  - The sentence (or question) and context length distributions are quite spread out, and  
28                  some questions/paragraphs are quite long. Classic Recurrent Neural Networks are  
29                  known to usually perform poorly on long sentences (it is observed in Machine  
30                  Translation for instance), because in practice they fail at keeping track of too much  
31                  information. To resolve this we will use Long Short Term Memory (LSTM) cells in

- 32 our model instead of classic RNN cells.  
 33 - The answers are often short. To discuss more what are the implications let's look at  
 34 the types of questions we have in our dataset

35  
 36 **2.1 Questions type**  
 37

Question	Count	Rate
How?	7545	9 %
Where?	3015	4 %
What?	34729	43 %
Who?	7479	9 %
Why?	1102	1 %

38 Interestingly, most of our train-set's questions are what-type questions. We thus expected to  
 39 perform quite well on this type of questions on the validation or test set.

40 Besides, it is a fact that where/what/who-type questions generally lead to shorter answers  
 41 compared to why-type questions. It makes sense that we observed a spike on short answers-  
 42 length in our previous graph. We thus expect our model to perform better on short answers  
 43 than long answers. As most of answers are short, we will try to find a trick so that our model  
 44 learns to prioritize some specific parts of the input and focus more on these parts.

45  
 46 **1.3 Important details to keep track of**  
 47

During our study, we will constantly keep track of the following behaviors in our model:

- Gradient saturation: We will pay extra attention when initializing the weights of sigmoid/tanh neurons to prevent saturation. Besides, we will constantly print the gradient at the sigmoid/tanh neurons and keep track of its values.
- Exploding gradient: we will use gradient clipping (by ceiling the gradient value at 10) to avoid an exploding gradient to compromise the convergence of our model.
- Size of the hidden state in our LSTM cell: It is very important to test models with different hidden state sizes for the LSTM cell. Indeed, increasing the state size gives more flexibility to the model and better learning capacity (because we have more parameters). However, it increases the complexity of the model, the number of computations grows exponentially and it becomes hard to train the model in a finite amount of time. We have to find the right tradeoff.
- Number of stacked LSTM: Just as we said for the size of the hidden state in our LSTM cell, we will have to find a tradeoff between accuracy and complexity of the model by tuning the number of stacked LSTM layers.
- Dropout probability: We will use dropout and we will train our model with different dropout probabilities to make sure we are not overfitting the data.
- Overfitting the data: We will repeatedly train the model on a small sample dataset and try to overfit the data, to make sure our model has enough capacity to understand complex behavior. This doesn't mean the model is good but it gives us an idea on its flexibility.

68  
 69 **2 Baseline model**

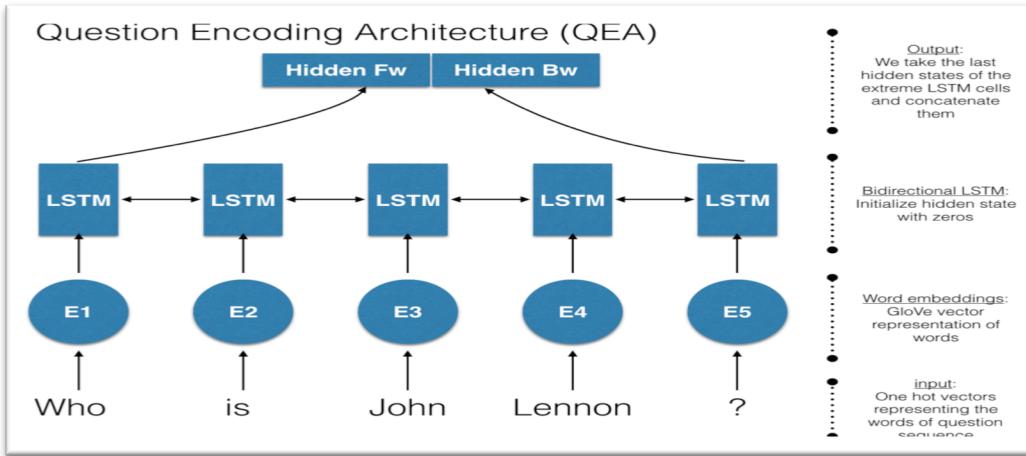
70 It is now time to jump on the detail of our first model. Our goal is to:

- Build a simple model that runs correctly in tensorflow. Once we get such a model, we can explore further and build more complicated features upon it.
- Evaluate this simple model with specific metrics, namely F1 and EM scores

74  
 75 **2.1 Question Encoding Architecture (QEA)**

76 To understand our Baseline model, we need to understand each of its

77 components, starting with the Question Encoding Architecture (QEA).



78

79 The input of the QEA is a sequence of one-hot vectors corresponding to the words of a given  
80 question. We first convert these words in their GloVe representation and give that as input to  
81 a Bidirectional LSTM where the hidden states are initialized as zero-vectors.

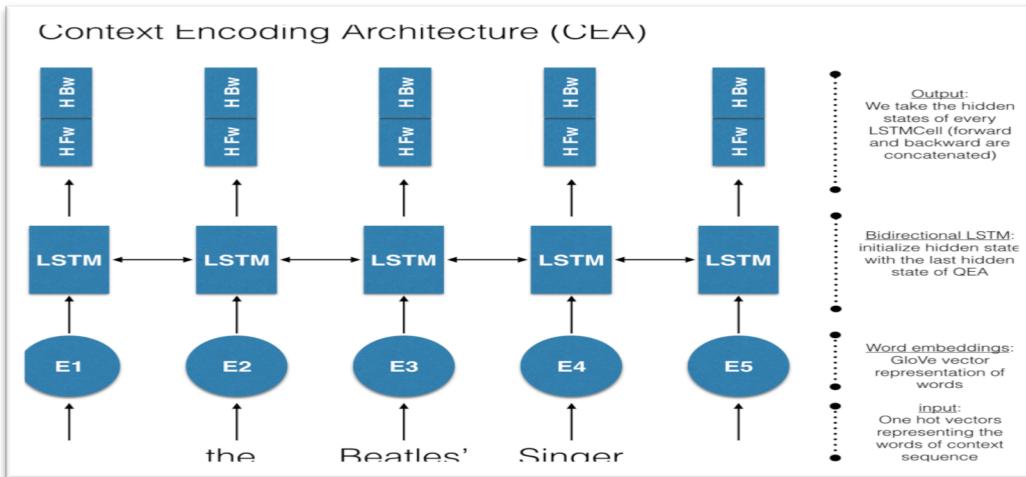
82 We compute the output of the QEA by concatenating the hidden states of the extremities of  
83 the Bi-LSTM. This output (output\_QEA) will be used two times:

- 84 - to initialize the hidden states of the Bi-LSTM of the context words  
85 - as an input to the decoder

86

## 87 2.2 Context Encoding Architecture (CEA)

88 As we said, we have another Bi-LSTM based model that takes the context words (GloVe  
89 Vectors) as inputs and has its hidden states initialized using the output\_QEA.



90

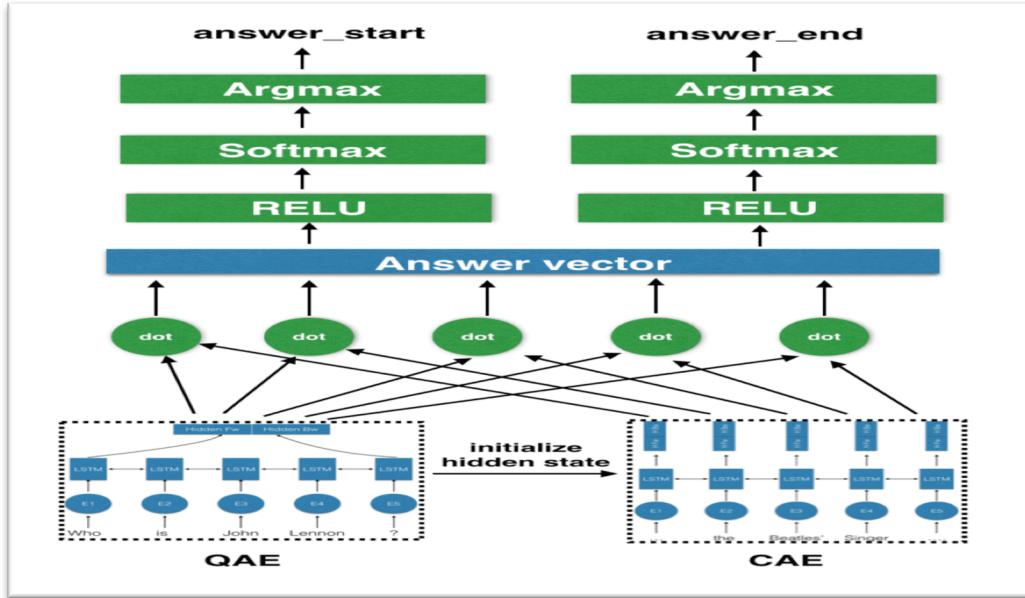
91 The CEA is very similar to the QAE, apart from the fact that we record the hidden-state  
92 vector of every LSTM cell.

93

## 94 2.3 Baseline Model Architecture (BMA)

95 After encoding our question sequence (QEA) and our context paragraph (CEA), we will use  
96 these new representation to find the answer. To do that, we will perform several layers of

97 computations as explained in the following scheme.



98

99 To keep track of the information coming from both the question and the context, we will  
100 first perform a dot product between QEA output and each CEA output. This is how we  
101 extract a feature answer vector. On this vector we successively apply a fully-connected layer  
102 with ReLU, Softmax and finally taking the argmax to find the indexes answer\_start and  
103 answer\_end.

104

## 105 2.4 Important details

106 In this Baseline Model Architecture, we also looked at the following issues and tricks:

- 107 - Exponential decay of the learning rate: We wanted the learning rate to be fast at the  
108 beginning of the training, because we start from random, we rapidly want get to a  
109 smart model and then optimize it more and more slowly.
- 110 - Masks: We want to keep the right vector/matrix size for our representation given  
111 any (question, context) pair. Without a mask, our LSTM trains on some random  
112 input figures and the loss is wrongly evaluated.
- 113 - Dropout layers: to avoid overfitting we added a dropout layer in the fully connected  
114 layer with ReLU.
- 115 - Gradient clipping: We ceiled the gradients in all the backpropagation so that there is  
116 no exploding gradients that would fail our learning.
- 117 - Batch size tuning: The higher the batch size, the lower the training time and the  
118 higher the memory needed. We tuned this batch size until we found the right  
119 tradeoff given the time/memory constraints.

120

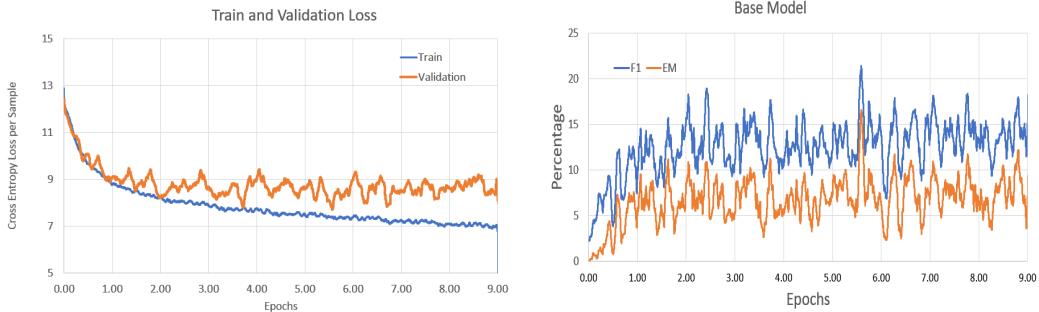
## 121 2.5 Baseline model results

122 Let's run our baseline model on the training set and observe the learning as well as the  
123 predictions.

124

125

126



127

128 We clearly see a good learning capability on the Train and Validation cross-entropy loss  
 129 plot. The learning capability echoes the second plot as well with a growing F1 score and a  
 130 oscillation around 14%.

131 Having a working baseline model makes the next steps a lot easier, and we only have to  
 132 work on the model parts (encoding + decoding). Let's now try to improve this model.

133

### 134 **3 Baseline model explorations**

135 In this part, we are going to underline some explorations that we did.

136

#### 137 **3.1 Stacked-LSTM Depth**

138 As said in part (1.3), adding several layers of Bi-LSTM one after each other makes our  
 139 model more complex and capable of extracting more information. However, it also makes  
 140 the time complexity increase. Let's look at our analysis for different number of Bi-LSTM  
 141 layers.

#Stacks	1	2	3
Performance	Correct	Better	overcomplexity in time

142 We decided to keep 1 Stack because of the time constraint. Even if the results were better  
 143 using 2 Stacks, it was taking about twice the time to train our model on an epoch for 1 Stack.  
 144 See plots in (3.3).

145

#### 146 **3.2 Hidden state size of the LSTM cell**

147 As said in part (1.3), it is very important to choose accurately the state size for the LSTM  
 148 cells as it can change the performances. However, it has the negative effect to drain a lot of  
 149 memory. Here's our observation:

150

151

152

153

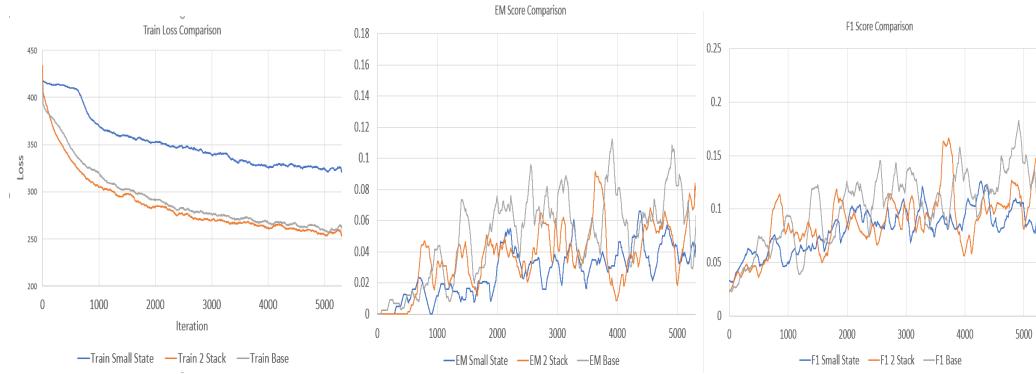
State size	64	128	256
Performance	Low	Correct	Out of memory

154 We decided to move to a state size of 128 for our LSTM cell because we still had enough  
 155 memory and were getting better results as you can see on the plot in (3.3). For 256, we  
 156 unfortunately ran out of memory as the number of parameters in our model increased  
 157 exponentially.

158

#### 159 **3.3 Exploration plots**

160



161

162

163 We can compare the performances between the baseline model (#State=128, 1 Stack) to a  
 164 model with (#State=128, 2 Stack) and another with (#State=64, 1 Stack). This is one of the  
 165 graph we used to tune these hyperparameters.

166

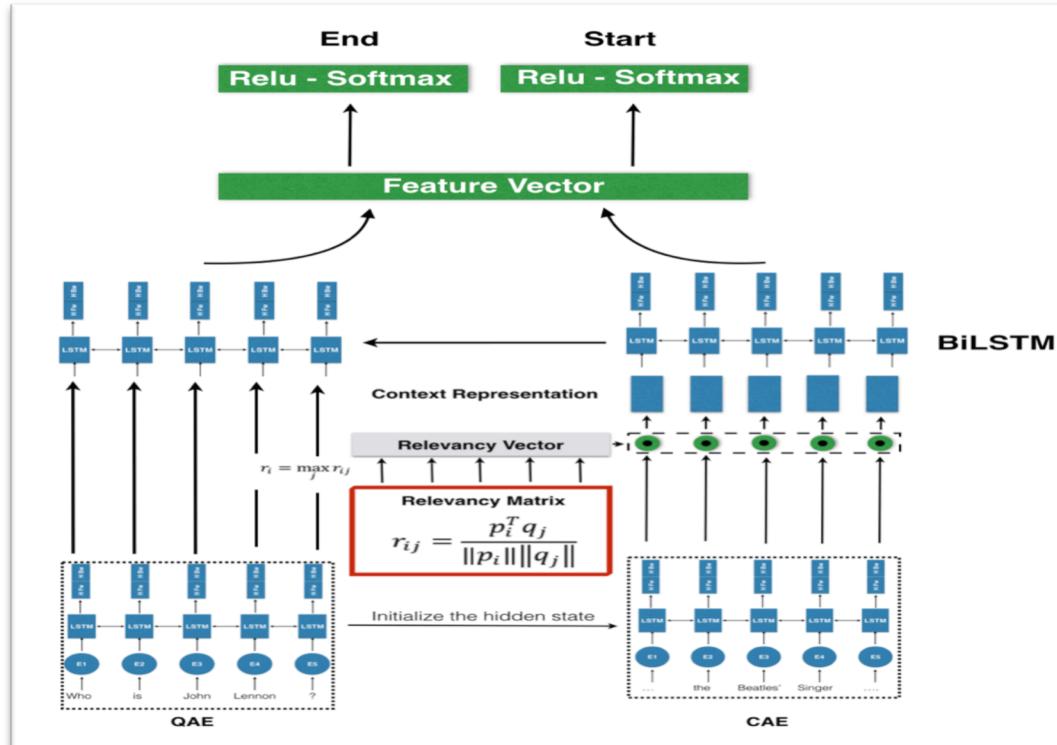
## 167 4 Personalized attention model

168 Armed with the knowledge based on Baseline model training, we propose a new network to  
 169 predict answer of a given question our model is based a paper from IBM Research ([1]), with  
 170 some improvement that we think would make the model better.

171

### 172 4.1 Model description

173 The following figure shows the abstract version of our model, given a question  $\{q_1, \dots, q_m\}$   
 174 and a paragraph  $\{p_1, \dots, p_n\}$ , we want to predict two number  $a_s$  and  $a_e$  as the index of start  
 175 and end of the question.



176

177   **Word embedding Layer**, in this layer for each word in the question and paragraph we  
 178   choose their corresponding word vector from Glove data set. This layer use embedding look  
 179   up table to retrieve the vector.

180   **Summarization layer**, we believe that given a question, each word of the context can be  
 181   summarized in a better way, to capture its importance the context. So, we first initialize the  
 182   BiLSTM with state size 128, to zero and feed a word embedding of question to that, to that  
 183   we obtain:

$$\overrightarrow{h_i^q} = \overrightarrow{\text{LSTM}}(\overrightarrow{h_{i-1}^q}, q_i) \quad \text{and}, \quad \overleftarrow{h_i^q} = \overleftarrow{\text{LSTM}}(\overleftarrow{h_{i-1}^q}, q_i)$$

184   The same LSTM layer is applied to the paragraph initializing by the final backward and  
 185   forward hidden state of question, so we obtain the same for paragraph:

$$\overrightarrow{h_i^p} = \overrightarrow{\text{LSTM}}(\overrightarrow{h_{i-1}^p}, p_i) \quad \text{and}, \quad \overleftarrow{h_i^p} = \overleftarrow{\text{LSTM}}(\overleftarrow{h_{i-1}^p}, p_i)$$

186   One of our major contribution to the main paper, is using this summarization layer before  
 187   amplification layer.

188   **Amplification layer**, using the output of summarization layer, by concatenating the results  
 189   of backward and forward hidden state, we can find a new representation each question word  
 190   and paragraph word as  $\{q'_1, \dots, q'_m\}$  and  $\{p'_1, \dots, p'_n\}$ , then by using cosine similarity measure  
 191   we can identify which words in the context are more important, it's kind of variant to  
 192   attention:

$$r_{ij} = \frac{q_i'^T p_j'}{\|q_i'\| \|p_j'\|}$$

193   By this similarity matrix we will amplify each word of the context by  $r_j = \max_i r_{ij}$ , that is  
 194   we build new representation by

$$p_i^{amp} = r_i * p_i$$

195   **Abstraction Layer**, after obtaining new representation of context words, based on their  
 196   importance, we use another BiLSTM layer to encode the paragraph and question, this time  
 197   we use the other way of feeding, that is we first encode the paragraph, using  
 198    $\{p_1^{amp}, \dots, p_n^{amp}\}$ , then feed the final backward and forward hidden state to encode the  
 199   question with the same BiLSTM with hidden state size of 128, so we obtain:

$$\begin{aligned} \overrightarrow{h_i^{p*}} &= \overrightarrow{\text{LSTM}}(\overrightarrow{h_{i-1}^{p*}}, p_i^{amp}) \quad \text{and}, \quad \overleftarrow{h_i^{p*}} = \overleftarrow{\text{LSTM}}(\overleftarrow{h_{i-1}^{p*}}, p_i^{amp}) \\ \overrightarrow{h_i^q} &= \overrightarrow{\text{LSTM}}(\overrightarrow{h_{i-1}^q}, q_i) \quad \text{and}, \quad \overleftarrow{h_i^q} = \overleftarrow{\text{LSTM}}(\overleftarrow{h_{i-1}^q}, q_i) \end{aligned}$$

200   Now we can use this hidden state representation of each word (by adding them together,  
 201   namely  $P_i = \frac{1}{2}(\overrightarrow{h_i^{p*}} * \overleftarrow{h_i^{p*}})$ ) to obtain some feature for each word in paragraph.

202   **Feature Vector**, We define six different feature for each word of paragraph as follows:

$$\begin{aligned} m_1^j &= \frac{1}{M} \sum_i \overrightarrow{h_i^q}^T P_j, \quad m_2^j = \frac{1}{M} \sum_i \overleftarrow{h_i^q}^T P_j, \quad m_3^j = \max_i \overrightarrow{h_i^q}^T P_j \\ m_4^j &= \max_i \overleftarrow{h_i^q}^T P_j, \quad m_5^j = \overrightarrow{h_{end}}^T P_j, \quad m_6^j = \overleftarrow{h_{start}}^T P_j \end{aligned}$$

205   These 6 features for each word can be used to figure out whether a word is the start or end of an  
 206   answer.

207   **Classification**, to capture the dependency of neighbor words, we pass this feature vector to  
 208   another BiLSTM with state size 8, and we used the maximum of hidden state value to pass to fully  
 209   connected layer.

$$\begin{aligned} \overrightarrow{h_i^m} &= \overrightarrow{\text{LSTM}}(\overrightarrow{h_{i-1}^m}, m_i) \quad \text{and}, \quad \overleftarrow{h_i^m} = \overleftarrow{\text{LSTM}}(\overleftarrow{h_{i-1}^m}, m_i) \\ s_i &= \max H_i^m[j], \quad H_i^m = [\overrightarrow{h_i^m}, \overleftarrow{h_i^m}] \end{aligned}$$

210   **Classification**, Now, having a number  $s$  for each word of context we classify this using a

211 fully connected layer with ReLu activation:

$$p_1 = \text{ReLU}(sW_2 + b_2), p_2 = \text{ReLU}(sW_1 + b_1)$$

212 And then using a SoftMax layer to predict the probability of each word to be the start or the  
213 end of an answer.

$$a_s = \text{argmax}(p_1), a_e = \text{argmax}(p_2)$$

214

## 215 4.2 Results

216

217 After hours and hours of work on this model, we finally managed to get to  
218 an F1 score of 24% and EM score of 19%. But we are still stuck on some  
219 bugs that we weren't able to figure out yet. We are looking forward to big  
220 improvement on this model that adds an attention mechanism to our  
221 baseline.

222

223 Things that we noticed:

224

- It is important to define all the variable inside the qa scope.  
Otherwise, at each step we initialize the variables again.
- Dropout should take as argument "keep probability" instead of "drop probability".
- To make sure the LSTM will use the previous variables, we need to reuse the scope to share variables.

225

## 226 5 Conclusions

227

228 This study has led us to build a baseline model that predicts the answer in a context  
229 paragraph given a specific question. We have looked at many channel of improvement on  
230 this baseline model and ended up tuning several hyperparameters and catching bugs. We  
231 then undertook the built of another model, more complex, based on an attention mechanism.  
232 We tuned it as much as we could and observed results. We believe that with more time,  
233 computing power and help we can get very good results with our second model.

234

235

236

237

238

239

240

241

242

243

244

## 245 References

246 [1] Wang, Z., Mi, H., Hamza, W., & Florian, R. (2016). Multi-Perspective Context Matching for Machine  
247 Comprehension. *arXiv preprint arXiv:1612.04211*.

248

249 [2] Seo, M., Kembhavi, A., Farhadi, A., & Hajishirzi, H. (2016). Bidirectional Attention Flow for Machine  
Comprehension. *arXiv preprint arXiv:1611.01603*.

250

251 [3] Xiong, C., Zhong, V., & Socher, R. (2016). Dynamic Coattention Networks For Question Answering.  
*arXiv preprint arXiv:1611.01604*.