

Lecture 11: Question Answering

1. What is question answering?

The goal of question answering

- to build systems that automatically answer questions posed by humans in a natural language
- information source: a text passage, all web documents, knowledge bases, tables, images.
- question type: factoid vs non-factoid, open-domain vs closed-domain, simple vs compositional
- answer type: a short segment of text, a paragraph, a list, yes

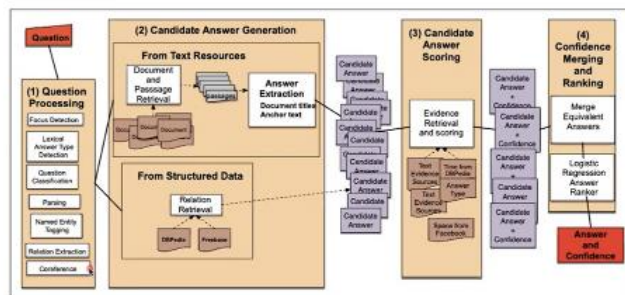


Image credit: J & M, edition 3

(1) Question processing, (2) Candidate answer generation, (3) Candidate answer scoring, and (4) Confidence merging and ranking.

2. Reading comprehension

Reading comprehension: building systems to comprehend a passage of text and answer questions about its content (P, Q)→A

Stanford question answering dataset (SQuAD)

- 100k annotated (passage, question, answer) triples
- Passages are selected from English Wikipedia, usually 100~150 words
- Questions are crowd-sourced
- Each answer is a short segment of text (or span) in the passage. → This is a limitation not all the questions can be answered in this way

Evaluation : exact match (0 or 1) and F1 (partial credit)

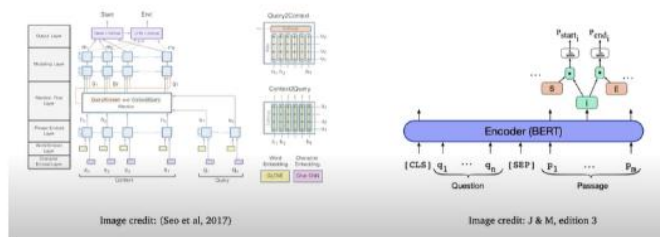
For development and testing sets, 3 gold answers are collected, because there could be multiple plausible answers

We compare the predicted answer to each gold answer and take max scores. Finally, we take the average of all the examples for both exact match and F1.

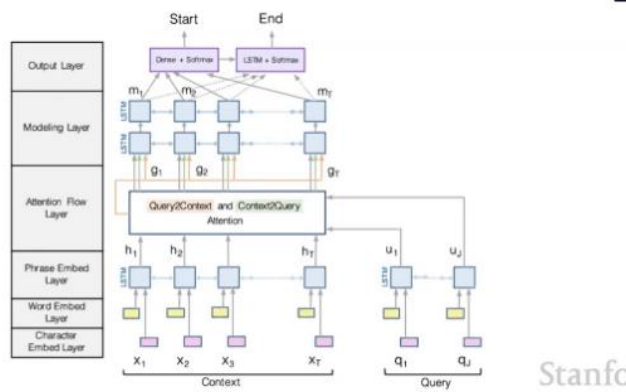
Build a model to solve SQuAD

- Problem formulation
 - input = $C = (c_1, c_2, c_3, \dots, c_N) = (q_1, q_2, \dots, q_M), c_i, q_i$
 - output: $1 \leq \text{start} \leq \text{end} \leq N$
 - $N \sim 100, M \sim 15$

LSTM-based vs BERT models

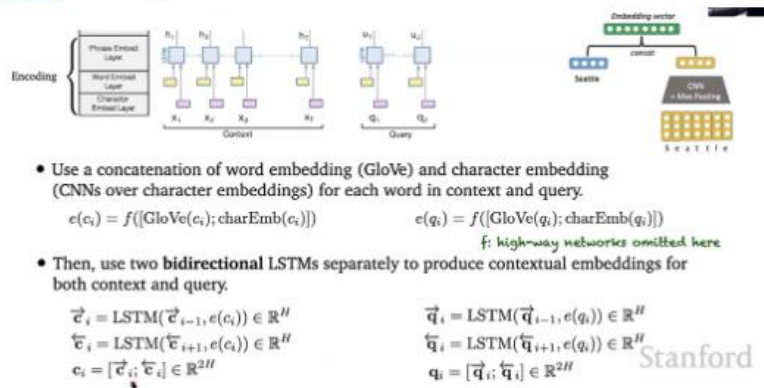


BiDAF: the Bidirectional Attention Flow model

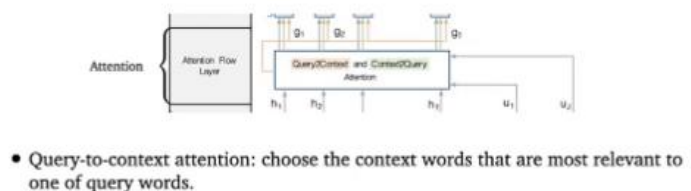


(Seo et al., 2017): Bidirectional Attention Flow for Machine Comprehension

BiDAF: Encoding



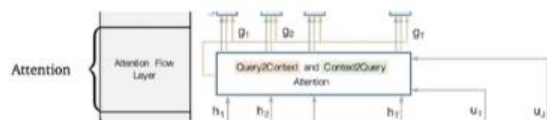
BiDAF: Attention



While Seattle's weather is very nice in summer, its weather is very rainy in winter making it one of the most gloomy cities in the U.S. LA is ...

Q: Which city is gloomy in winter?

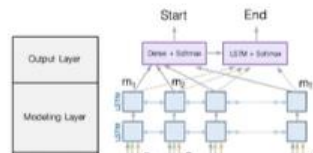
Stanf



- First, compute a similarity score for every pair of (c_i, q_j) :

$$S_{i,j} = \mathbf{w}_{\text{sim}}^T [\mathbf{c}_i; \mathbf{q}_j; \mathbf{c}_i \odot \mathbf{q}_j] \in \mathbb{R} \quad \mathbf{w}_{\text{sim}} \in \mathbb{R}^{6H}$$

BiDAF: Modeling and output layer

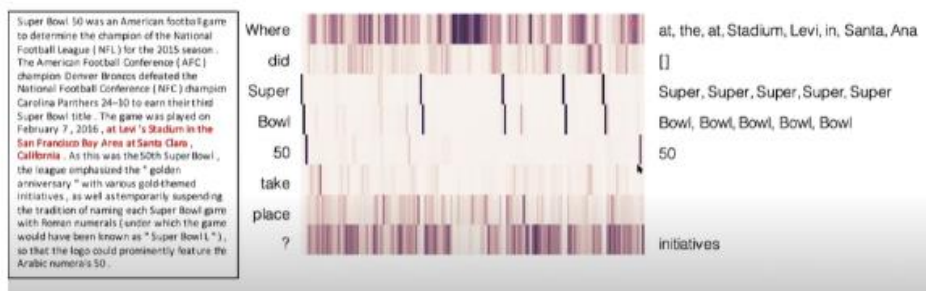


Modeling layer: pass \mathbf{g}_i to another two layers of bi-directional LSTMs.

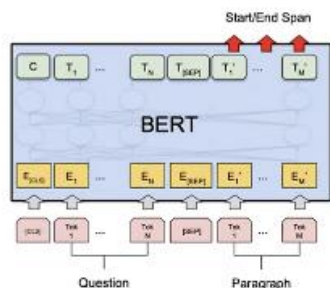
- Attention layer is modeling interactions between query and context
- Modeling layer is modeling interactions within context words

$$\mathbf{m}_i = \text{BiLSTM}(\mathbf{g}_i) \in \mathbb{R}^{2H}$$

attention visualization



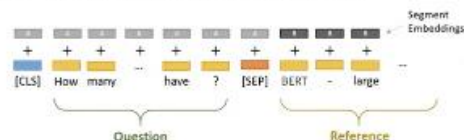
BERT for reading comprehension



Question = Segment A

Passage = Segment B

Answer = predicting two endpoints in segment B



Question: How many parameters does BERT-large have?

Reference Text: BERT-large is really big... it has 24 layers and an embedding size of 1,024, for a total of 340M parameters! Altogether it is 1.34GB, so expect it to take a couple minutes to download to your Colab instance.

Image credit: <https://mccormickml.com/>

$$\mathcal{L} = -\log p_{\text{start}}(s^*) - \log p_{\text{end}}(e^*)$$

$$p_{\text{start}}(i) = \text{softmax}_i(\mathbf{w}_{\text{start}}^T \mathbf{h}_i)$$

$$p_{\text{end}}(i) = \text{softmax}_i(\mathbf{w}_{\text{end}}^T \mathbf{h}_i)$$

where \mathbf{h}_i is the hidden vector of c_p returned by BERT

3. Open-domain (textual) question answering

- Different from reading comprehension, we don't assume a given passage.
- Instead, we only have access to a large collection of documents. We don't know where the answer is located, and the goal is to return the answer for any open-domain questions.
- Much more challenging but a more practical problem

Retriever-reader framework

