# Textual Question Answering

## SQuAD 1.1

输入T和Q，输出回答A。A必须是T中一个span，即extractive question answering。

### SQuAD evaluation, v1.1

- Authors collected 3 gold answers
- Systems are scored on two metrics:
  - Exact match: 1/0 accuracy on whether you match one of the 3 answers
  - F1: Take system and each gold answer as bag of words, evaluate
    Precision = $\frac{TP}{TP+FP}$, Recall = $\frac{TP}{TP+FN}$, harmonic mean F1 = $\frac{2PR}{P+R}$
    Score is (macro-)average of per-question F1 scores
- F1 measure is seen as more reliable and taken as primary
  - It's less based on choosing exactly the same span that humans chose, which is susceptible to various effects, including line breaks
- Both metrics ignore punctuation and articles (**a, an, the** only)

**缺点**

- 所有问题都有答案
- 系统只需要对span排序，并取分数最高的

## SQuAD 2.0

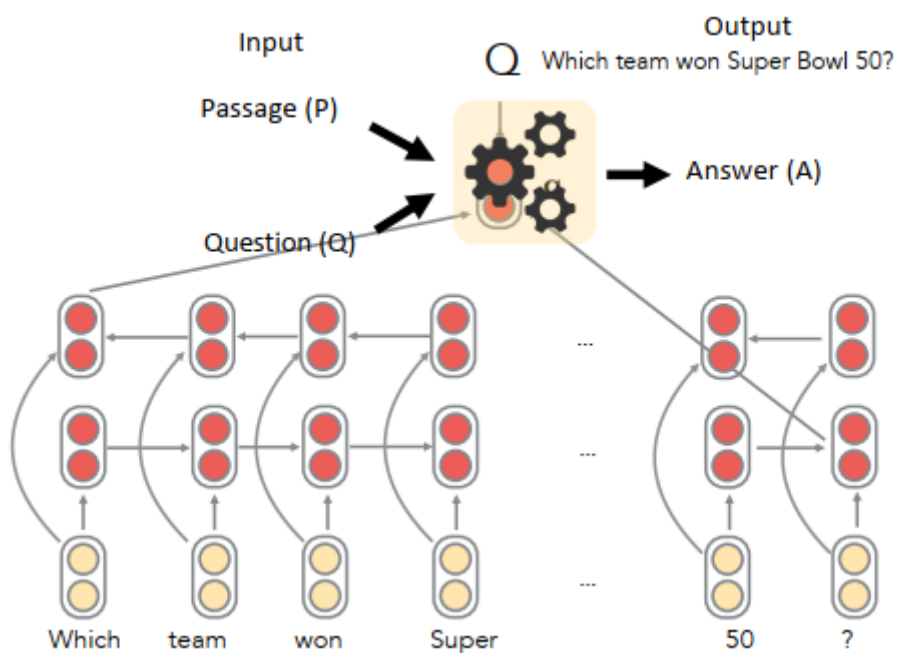train中1/3是no answer；dev/test中1/2是no answer。

**缺点**

- 答案必须是一段span，不能yes/no、计数等
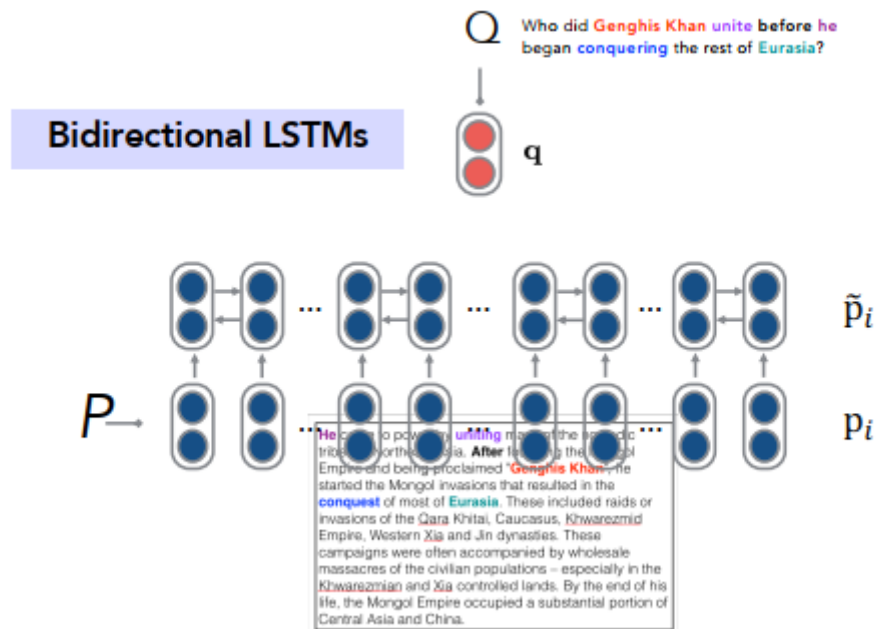- 不能有隐含意义，像"作者xxx的举动反应他的什么情绪"

但SQuAD仍是最广泛使用的QA数据集。
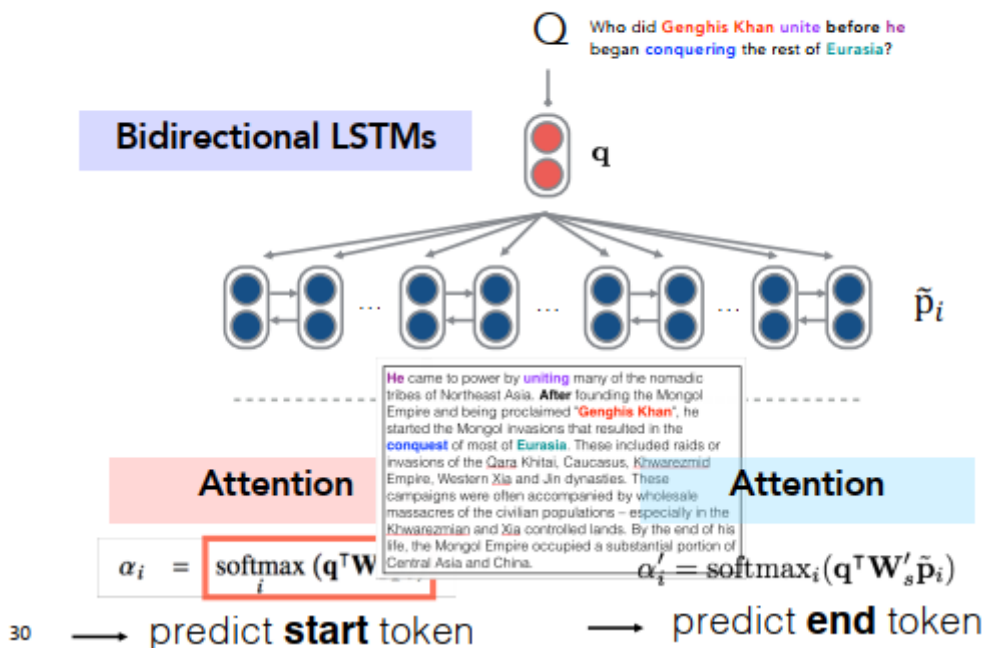
# Stanford Attentive Reader

具体算法流程:

1. 通过bi-lstm计算Q的向量:



2. 通过另一个bi-lstm计算每个P中单词的向量:

3. Q向量和P中单词向量计算attention分数，并得到A的<START>和<END>位置:



**Stanford Attentive Reader**

$$\alpha_i = \text{softmax}_i(\mathbf{q}^\top \mathbf{W} \ldots)$$

$$\alpha_i' = \text{softmax}_i(\mathbf{q}^\top \mathbf{W}_s' \tilde{\mathbf{p}}_i)$$

30 ⟶ predict **start** token ⟶ predict **end** token

# Stanford Attentive Reader++

2点改进:

- SAR用正向lstm和反向lstm的最后一层的hidden state作为Q向量。SAR++用bi-lstm每一层hidden state的加权和作为Q向量。
- wordvec加入manual特征，包括POS tags/ frequency/exact match等

# BiDAF (Bi-Directional Attention Flow)

## Attention Flow

之前各个context word对query的attention，而attention可以是双向流动的。BiDAF考虑context2query和query2context的attention。

## Dynamic Coattention Networks

C2Q+Q2C+second-level attention

## FusionNet

multi-level attention