

Deep Neural Networks for Natural Language Processing (AI6127)

JUNG-JAE KIM

LECTURE 9: (TEXTUAL) QUESTION ANSWERING

Lecture Plan

- Motivation/History
- The SQuAD dataset
- The Stanford Attentive Reader model
- BiDAF
- Recent, more advanced architectures
- ELMo and BERT preview



About 6,030,000 results (0.69 seconds)

John Christian Watson

John Christian Watson (born **John Christian Tanck**; 9 April 1867 – 18 November 1941), commonly known as **Chris Watson**, was an Australian politician who served as the third Prime Minister of Australia.



en.wikipedia.org

[Chris Watson - Wikipedia](https://en.wikipedia.org/wiki/Chris_Watson)https://en.wikipedia.org/wiki/Chris_Watson

People also search for

[View 15+ more](#)

Andrew Fisher



George Reid



Billy Hughes



Edmund Barton



Alfred Deakin



Kevin Rudd



Julia Gillard

[More about Chris Watson](#)

1. Motivation: Question answering

- With massive collections of full-text documents, i.e., the web, simply returning relevant documents is of limited use
- Rather, we often want **answers** to our **questions**
- Especially on mobile
- Or using a digital assistant device, like Alexa, Google Assistant, ...
- We can factor this into two parts:
 - Finding documents that (might) contain an answer
 - Which can be handled by traditional information retrieval/web search
 - Finding an answer in a paragraph or a document
 - This problem is often termed **Reading Comprehension**. It is what we will focus on today

A Brief History of Reading Comprehension

- Much early NLP work attempted reading comprehension
 - Schank, Abelson, Lehnert et al. c. 1977 – “Yale A.I. Project”
- Revived by Lynette Hirschman in 1999:
 - Could NLP systems answer human reading comprehension questions for 3rd to 6th graders? Simple methods attempted.
- Revived again by Chris Burges in 2013 with MCTest
 - Again answering questions over simple story texts
- Floodgates opened in 2015/16 with the production of large datasets which permit supervised neural systems to be built
 - Hermann et al. (NIPS 2015) DeepMind CNN/DM dataset
 - Rajpurkar et al. (EMNLP 2016) SQuAD
 - MS MARCO, TriviaQA, RACE, NewsQA, NarrativeQA, ...

Machine Comprehension (Burges 2013)

- “A machine comprehends a passage of text if, for any question regarding that text that can be answered correctly by a majority of native speakers, that machine can provide a string which those speakers would agree both answers that question, and does not contain information irrelevant to that question.”

Towards the Machine Comprehension of Text: An Essay

Christopher J.C. Burges
Microsoft Research
One Microsoft Way
Redmond, WA 98052, USA

December 23, 2013



MCTest Reading Comprehension

Passage (*P*) + Question (*Q*) → Answer (*A*)

P

Alyssa got to the beach after a long trip. She's from Charlotte. She traveled from Atlanta. She's now in Miami. She went to Miami to visit some friends. But she wanted some time to herself at the beach, so she went there first. After going swimming and laying out, she went to her friend Ellen's house. Ellen greeted Alyssa and they both had some lemonade to drink. Alyssa called her friends Kristin and Rachel to meet at Ellen's house.....

Q

Why did Alyssa go to Miami?

A

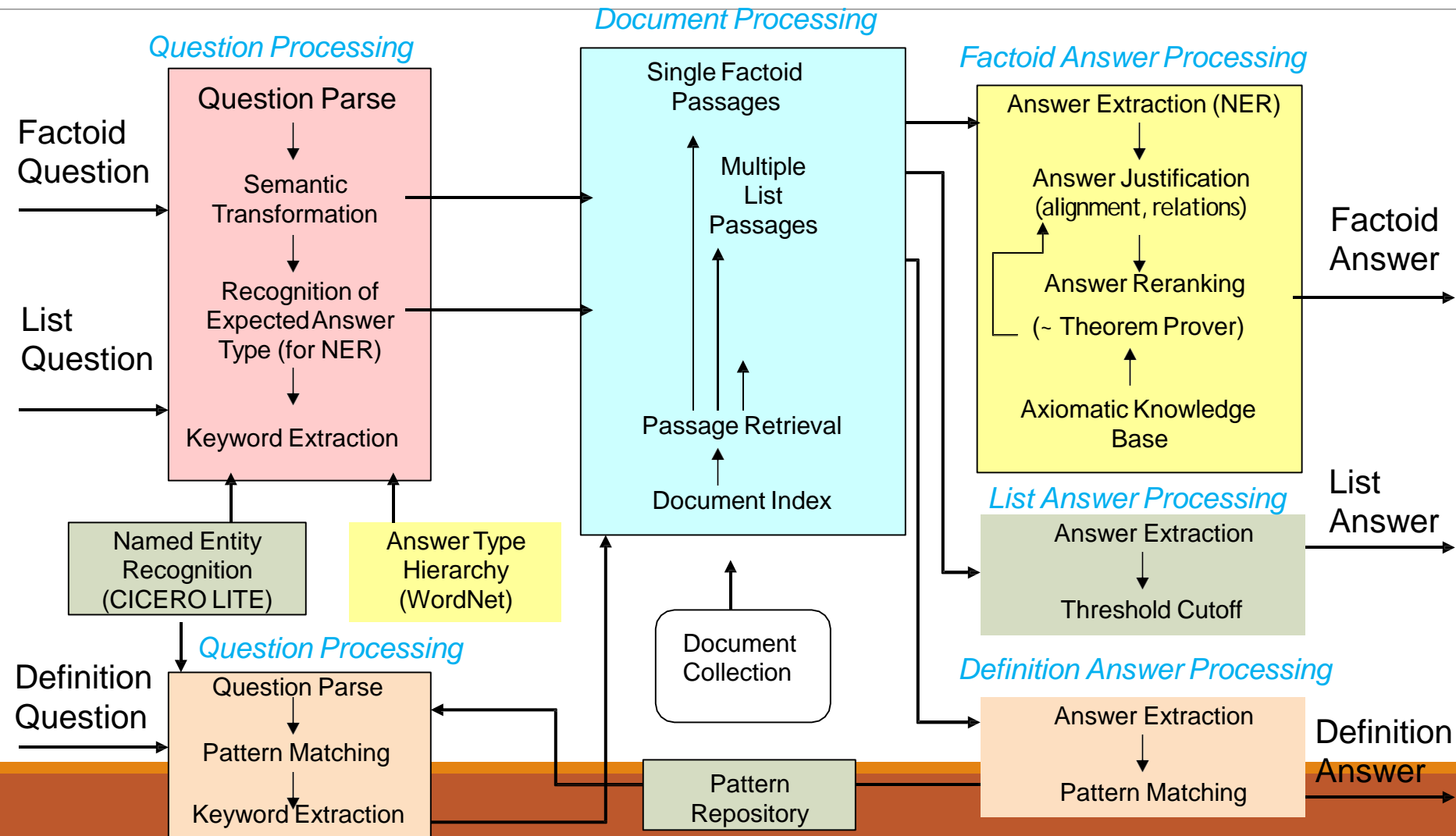
To visit some friends

A Brief History of Open-domain Question Answering

- Simmons et al. (1964) did first exploration of answering questions from an expository text based on matching dependency parses of a question and answer
- Murax (Kupiec 1993) aimed to answer questions over an online encyclopedia using IR and shallow linguistic processing
- The NIST TREC QA track begun in 1999 first rigorously investigated answering fact questions over a large collection of documents
- IBM's Jeopardy! System (DeepQA, 2011) brought attention to a version of the problem; it used an ensemble of many methods
- DrQA (Chen et al. 2016) used IR followed by neural reading comprehension to bring deep learning to Open-domain QA

[architecture of LCC (Harabagiu/Moldovan) QA system, circa 2003] Complex systems but they did work fairly well on “factoid” questions (e.g. who, what)

Turn-of-the Millennium Full NLP QA:



2. Stanford Question Answering Dataset (SQuAD)

(Rajpurkar et al., 2016)

Question: Which team won Super Bowl 50?

Passage

Super Bowl 50 was an American football game to determine the champion of the National Football League (NFL) for the 2015 season. The American Football Conference (AFC) champion **Denver Broncos** defeated the National Football Conference (NFC) champion Carolina Panthers 24–10 to earn their third Super Bowl title. The game was played on February 7, 2016, at Levi's Stadium in the San Francisco Bay Area at Santa Clara, California.

100k examples

Answer must be a span in the passage

A.k.a. extractive question answering

Stanford Question Answering Dataset (SQuAD)

Private schools, also known as independent schools, non-governmental, or nonstate schools, are not administered by local, state or national governments; thus, they retain the right to select their students and are funded in whole or in part by charging their students tuition, rather than relying on mandatory taxation through public (government) funding; at some private schools students may be able to get a scholarship, which makes the cost cheaper, depending on a talent the student may have (e.g. sport scholarship, art scholarship, academic scholarship), financial need, or tax credit scholarships that might be available.

Along with non-governmental and nonstate schools, what is another name for private schools?

Gold answers: (1) independent, (2) independent schools, (3) independent schools

Along with sport and art, what is a type of talent scholarship?

Gold answers: (1) academic, (2) academic, (3) academic

Rather than taxation, what are private schools largely funded by?

Gold answers: (1) tuition, (2) charging their students tuition, (3) tuition

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
$$precision = \frac{TP}{TP + FP} \quad recall = \frac{TP}{TP + FN} \quad F1 = \frac{2PR}{P + R} \quad (\text{harmonic mean})$$
 - 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)

SQuAD v1.1 leaderboard, end of 2016 (Dec 6)

		EM	F1
11	Fine-Grained Gating Carnegie Mellon University (Yang et al. '16)	62.5	73.3
12	Dynamic Chunk Reader IBM (Yu & Zhang et al. '16)	62.5	71.0
13	Match-LSTM with Ans-Ptr (Boundary) Singapore Management University (Wang & Jiang '16)	60.5	70.7
14	Match-LSTM with Ans-Ptr (Sequence) Singapore Management University (Wang & Jiang '16)	54.5	67.7
15	Logistic Regression Baseline Stanford University (Rajpurkar et al. '16)	40.4	51.0
Will your model outperform humans on the QA task?			
	Human Performance Stanford University (Rajpurkar et al. '16)	82.3	91.2

SQuAD v1.1 leaderboard, 2019-02-07 – it's solved!

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar et al. '16)	82.304	91.221
1 Oct 05, 2018	BERT (ensemble) Google AI Language https://arxiv.org/abs/1810.04805	87.433	93.160
2 Oct 05, 2018	BERT (single model) Google AI Language https://arxiv.org/abs/1810.04805	85.083	91.835
2 Sep 09, 2018	nlnet (ensemble) Microsoft Research Asia	85.356	91.202
2 Sep 26, 2018	nlnet (ensemble) Microsoft Research Asia	85.954	91.677
3 Jul 11, 2018	QANet (ensemble) Google Brain & CMU	84.454	90.490
4 Jul 08, 2018	r-net (ensemble) Microsoft Research Asia	84.003	90.147
5 Mar 19, 2018	QANet (ensemble) Google Brain & CMU	83.877	89.737
5 Sep 09, 2018	nlnet (single model) Microsoft Research Asia	83.468	90.133

SQuAD v1.1 leaderboard, 2019-09-27 – actively being improved

Rank	Model	EM	F1
	Human Performance <i>Stanford University</i> (Rajpurkar et al. '16)	82.304	91.221
1 May 21, 2019	XLNet (single model) <i>Google Brain & CMU</i>	89.898	95.080
2 Aug 11, 2019	XLNET-123 (single model) <i>MST/EOI</i>	89.646	94.930
3 Jul 21, 2019	SpanBERT (single model) <i>FAIR & UW</i>	88.839	94.635
4 Jul 03, 2019	BERT+WWM+MT (single model) <i>Xiao Research</i>	88.650	94.393
5 Jul 21, 2019	Tuned BERT-1seq Large Cased (single model) <i>FAIR & UW</i>	87.465	93.294
6 Oct 05, 2018	BERT (ensemble) <i>Google AI Language</i> https://arxiv.org/abs/1810.04805	87.433	93.160

SQuAD 2.0

- A defect of SQuAD 1.0 is that all questions have an answer in the paragraph
- Systems (implicitly) rank candidates and choose the best one
- You don't have to judge whether a span answers the question
- In SQuAD 2.0, 1/3 of the training questions have no answer, and about 1/2 of the dev/test questions have no answer
 - For NoAnswer examples, NoAnswer receives a score of 1, and any other response gets 0, for both exact match and F1
- Simplest system approach to SQuAD 2.0:
 - Have a threshold score for whether a span answers a question
- Or you could have a second component that confirms answering
 - Like Natural Language Inference (NLI) or "Answer validation"

SQuAD 2.0 Example

Genghis Khan united the Mongol and Turkic tribes of the steppes and became Great Khan in 1206. He and his successors expanded the Mongol empire across Asia. Under the reign of Genghis' third son, Ögedei Khan, the Mongols destroyed the weakened Jin dynasty in 1234, conquering most of northern China. Ögedei offered his nephew Kublai a position in Xingzhou, Hebei. Kublai was unable to read Chinese but had several Han Chinese teachers attached to him since his early years by his mother Sorghaghtani. He sought the counsel of Chinese Buddhist and Confucian advisers. Möngke Khan succeeded Ögedei's son, Güyük, as Great Khan in 1251. He

When did Genghis Khan kill Great Khan?

Gold Answers: <No Answer>

Prediction: 1234 [from Microsoft nlnet]

SQuAD 2.0 leaderboard, BiDAF variants

36 Sep 13, 2018	BiDAF++ (single model) <i>UW and FAIR</i>	65.651	68.866
37 Jun 27, 2018	BSAE AddText (single model) <i>reciTAL.ai</i>	63.338	67.422
38 Aug 14, 2018	eeAttNet (single model) <i>BBD NLP Team</i> https://www.bbdservice.com	63.327	66.633
38 May 30, 2018	BiDAF + Self Attention + ELMo (single model) <i>Allen Institute for Artificial Intelligence</i> <i>[modified by Stanford]</i>	63.372	66.251
39 Nov 27, 2018	Tree-LSTM + BiDAF + ELMo (single model) <i>Carnegie Mellon University</i>	57.707	62.341
39 May 30, 2018	BiDAF + Self Attention (single model) <i>Allen Institute for Artificial Intelligence</i> <i>[modified by Stanford]</i>	59.332	62.305
40 May 30, 2018	BiDAF-No-Answer (single model) <i>University of Washington [modified by Stanford]</i>	59.174	62.093

<https://rajpurkar.github.io/SQuAD-explorer/>

SQuAD 2.0 leaderboard

2019-02-07

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452
1 Jan 15, 2019	BERT + MMFT + ADA (ensemble) Microsoft Research Asia	85.082	87.615
2 Jan 10, 2019	BERT + Synthetic Self-Training (ensemble) Google AI Language https://github.com/google-research/bert	84.292	86.967
3 Dec 13, 2018	BERT finetune baseline (ensemble) Anonymous	83.536	86.096
4 Dec 16, 2018	Lunet + Verifier + BERT (ensemble) Layer 6 AI NLP Team	83.469	86.043
4 Dec 21, 2018	PAML+BERT (ensemble model) PINGAN GammaLab	83.457	86.122
5 Dec 15, 2018	Lunet + Verifier + BERT (single model) Layer 6 AI NLP Team	82.995	86.035

SQuAD 2.0 leaderboard 2019-09-27, it's solved!

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452
1 Sep 18, 2019	ALBERT (ensemble model) Google Language ALBERT Team	89.731	92.215
2 Jul 22, 2019	XLNet + DAAF + Verifier (ensemble) PINGAN Omni-Sinitic	88.592	90.859
2 Sep 16, 2019	ALBERT (single model) Google Language ALBERT Team	88.107	90.902
2 Jul 26, 2019	UPM (ensemble) Anonymous	88.231	90.713
3 Aug 04, 2019	XLNet + SG-Net Verifier (ensemble) Shanghai Jiao Tong University & CloudWalk https://arxiv.org/abs/1908.05147	88.174	90.702
4 Aug 04, 2019	XLNet + SG-Net Verifier++ (single model) Shanghai Jiao Tong University & CloudWalk https://arxiv.org/abs/1908.05147	87.238	90.071
5 Jul 26, 2019	UPM (single model) Anonymous	87.193	89.934

Good systems are great, but still basic NLU errors

The Yuan dynasty is considered both a successor to the Mongol Empire and an imperial Chinese dynasty. It was the khanate ruled by the successors of Möngke Khan after the division of the Mongol Empire. In official Chinese histories, the Yuan dynasty bore the Mandate of Heaven, following the Song dynasty and preceding the Ming dynasty. The dynasty was established by Kublai Khan, yet he placed his grandfather Genghis Khan on the imperial records as the official founder of the

What dynasty came before the Yuan?

Gold Answers: (1) Song dynasty, (2) Mongol Empire,
(3) the Song dynasty

Prediction: Ming dynasty [BERT (single model) (GoogleAI)]

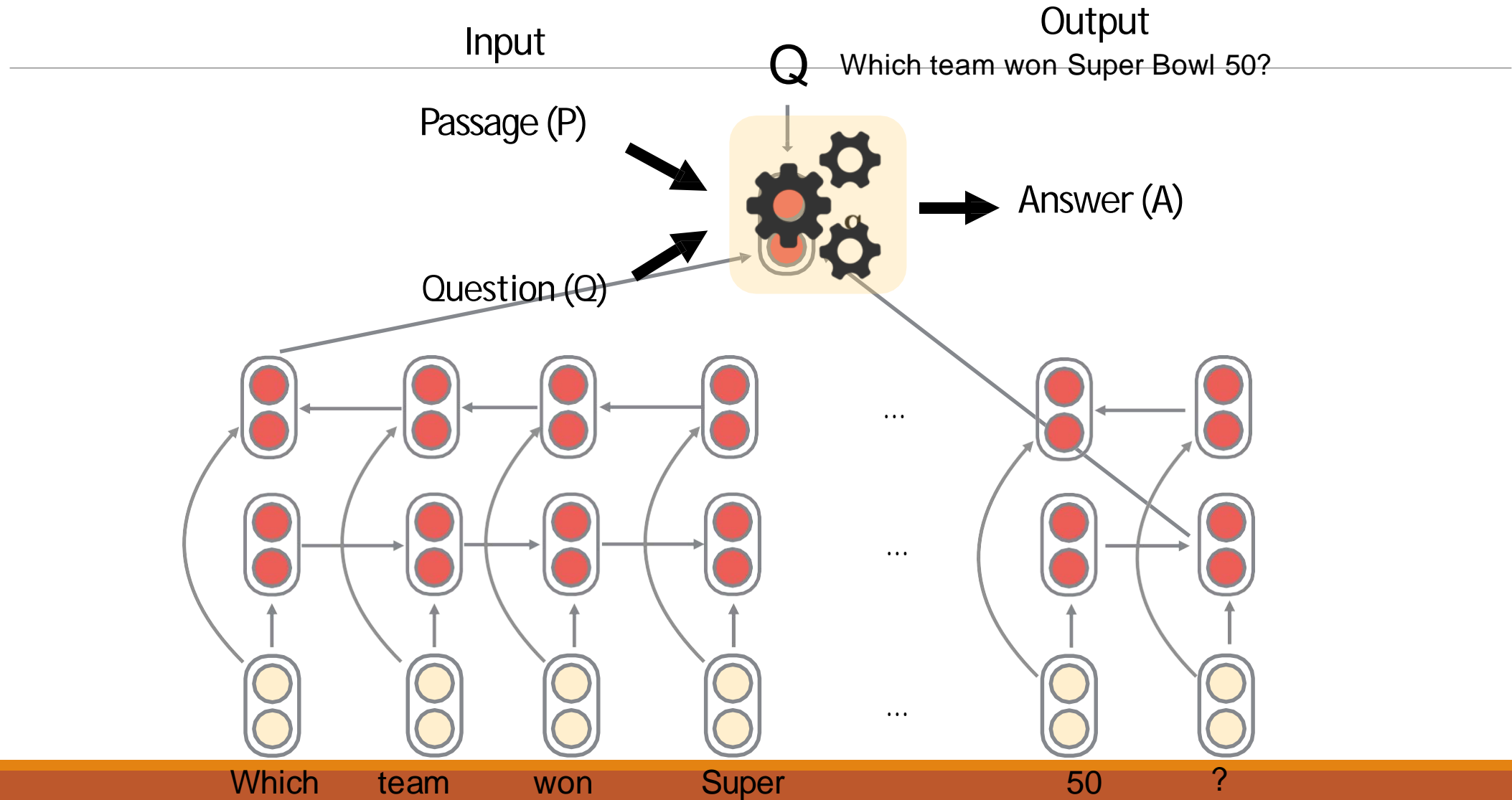
SQuAD limitations

- SQuAD has a number of other key limitations too:
 - Only span-based answers (no yes/no, counting, implicit why)
 - Questions were constructed looking at the passages
 - Not genuine information needs
 - Generally greater lexical and syntactic matching between questions and answer span than you get in real life
 - Barely any multi-fact/sentence inference beyond coreference
- Nevertheless, it is a well-targeted, well-structured, clean dataset
 - It has been the most used and competed-on QA dataset
 - It has also been a useful starting point for building systems in industry (though in-domain data always really helps!)

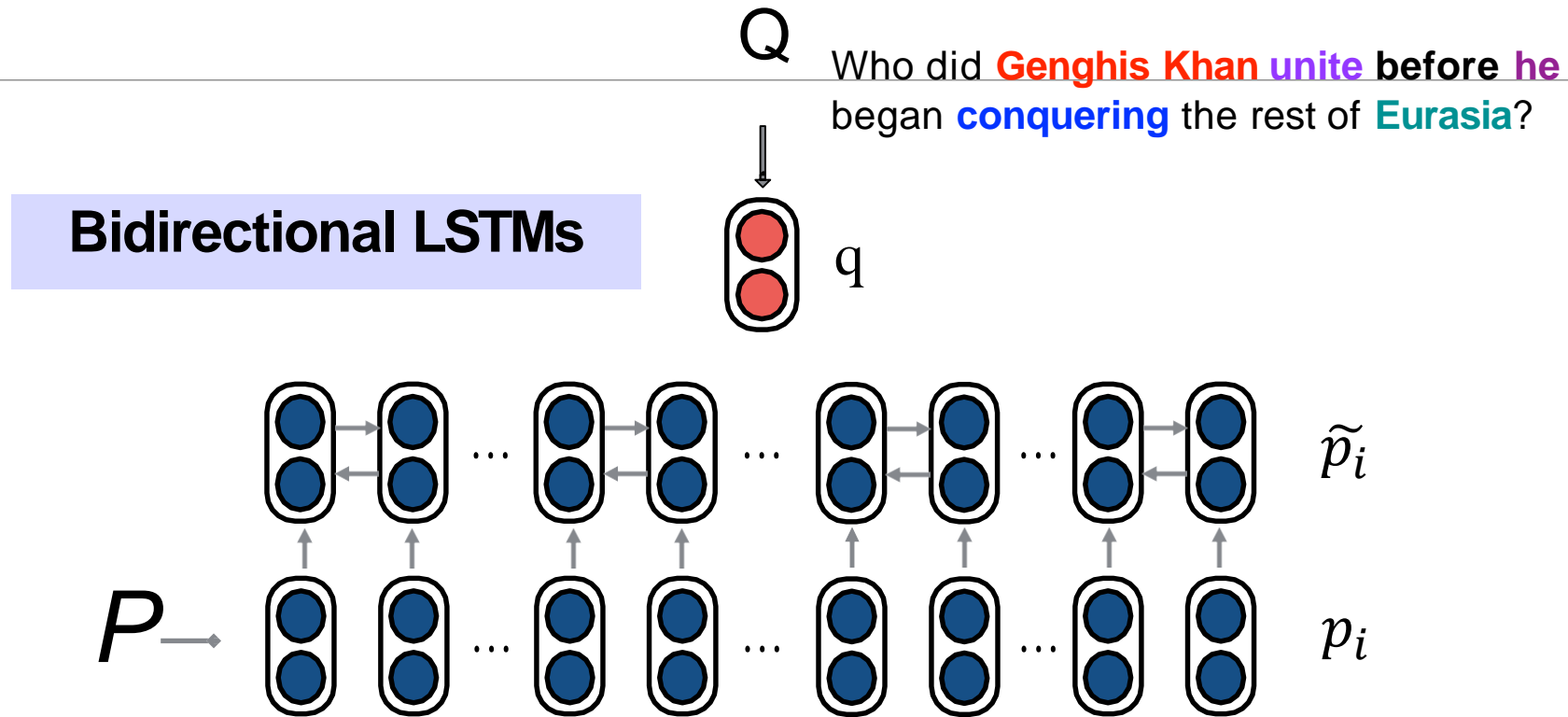
3. Stanford Attentive Reader

- Demonstrated a minimal, highly successful architecture for reading comprehension and question answering

The Stanford Attentive Reader

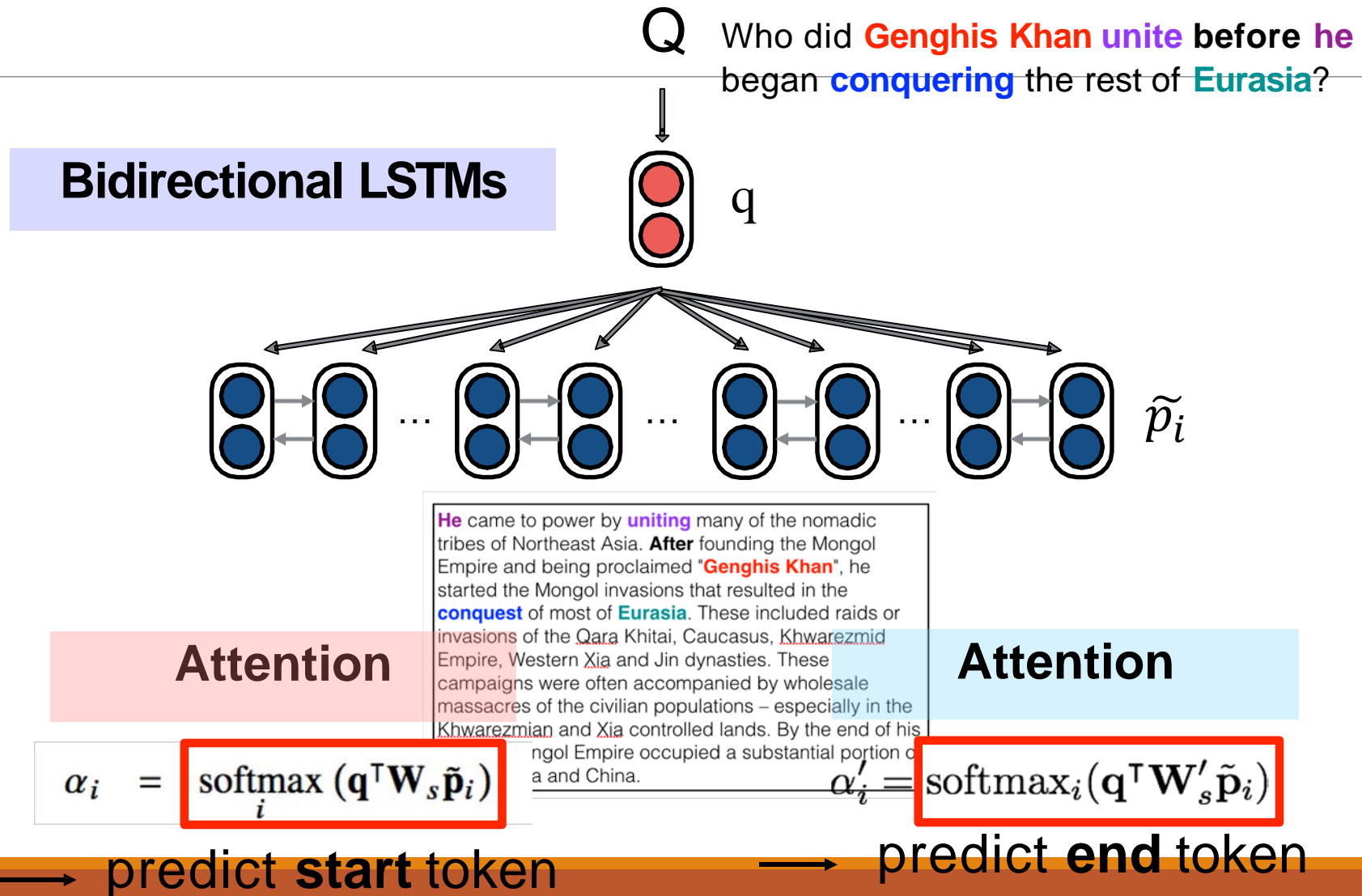


Stanford Attentive Reader



He came to power by **uniting** many of the nomadic tribes of Northeast Asia. **After** founding the Mongol Empire and being proclaimed "**Genghis Khan**", he started the Mongol invasions that resulted in the **conquest** of most of **Eurasia**. These included raids or invasions of the Qara Khitai, Caucasus, Khwarezmid Empire, Western Xia and Jin dynasties. These campaigns were often accompanied by wholesale massacres of the civilian populations – especially in the **Khwarezmian** and **Xia** controlled lands. By the end of his life, the Mongol Empire occupied a substantial portion of Central Asia and China.

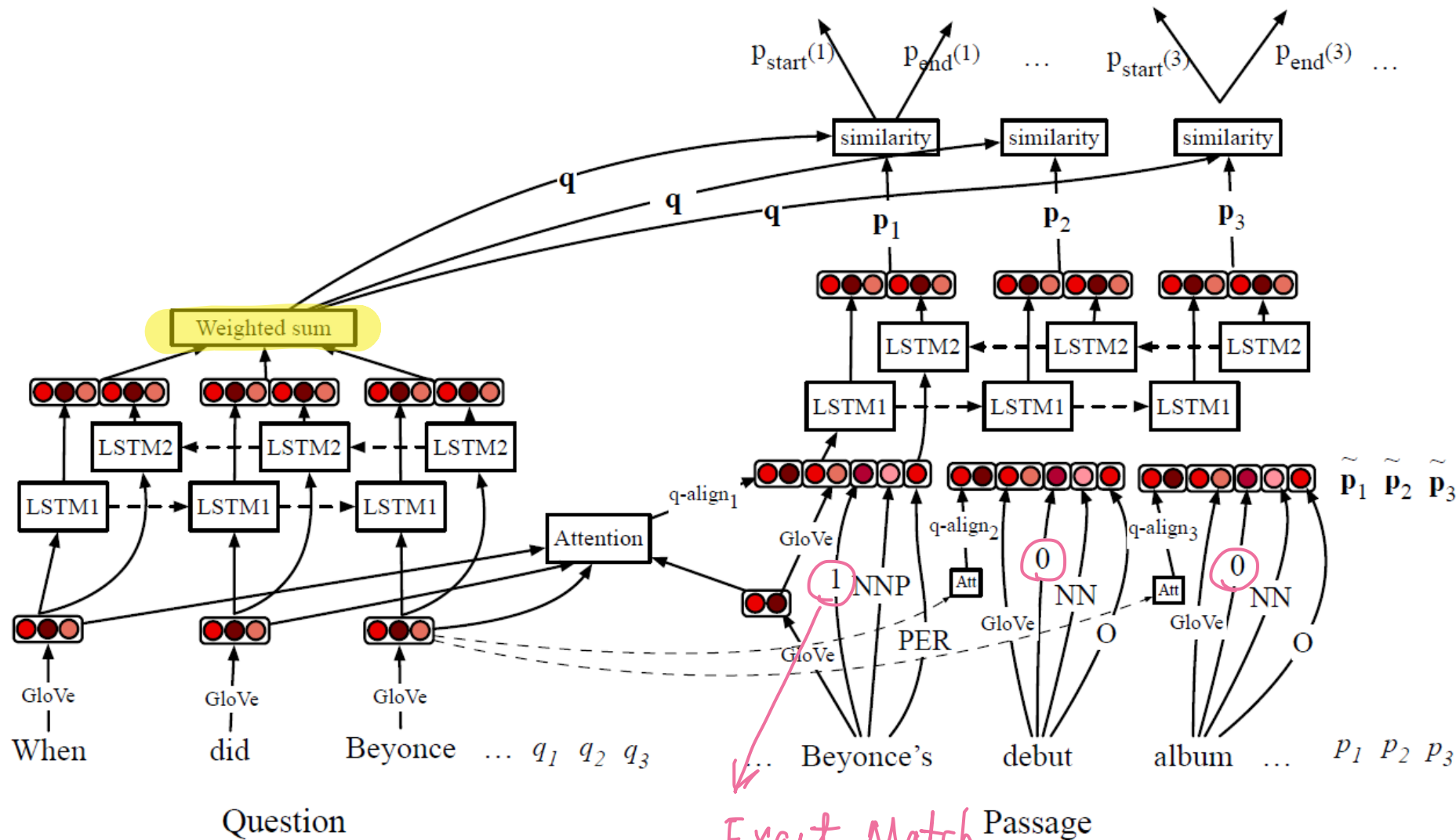
Stanford Attentive Reader



SQuAD 1.1 Results (single model, c. Feb 2017)

	F1
Logistic regression	51.0
Fine-Grained Gating (Carnegie Mellon U)	73.3
Match-LSTM (Singapore Management U)	73.7
DCN (Salesforce)	75.9
BiDAF (UW & Allen Institute)	77.3
Multi-Perspective Matching (IBM)	78.7
ReasoNet (MSR Redmond)	79.4
DrQA (Chen et al. 2017)	79.4
r-net (MSR Asia) [Wang et al., ACL2017]	79.7
Human performance	91.2

Stanford Attentive Reader++



- Word embedding
- Exact match (lemma) to question word
- POS tag
- NER tag

Training objective:

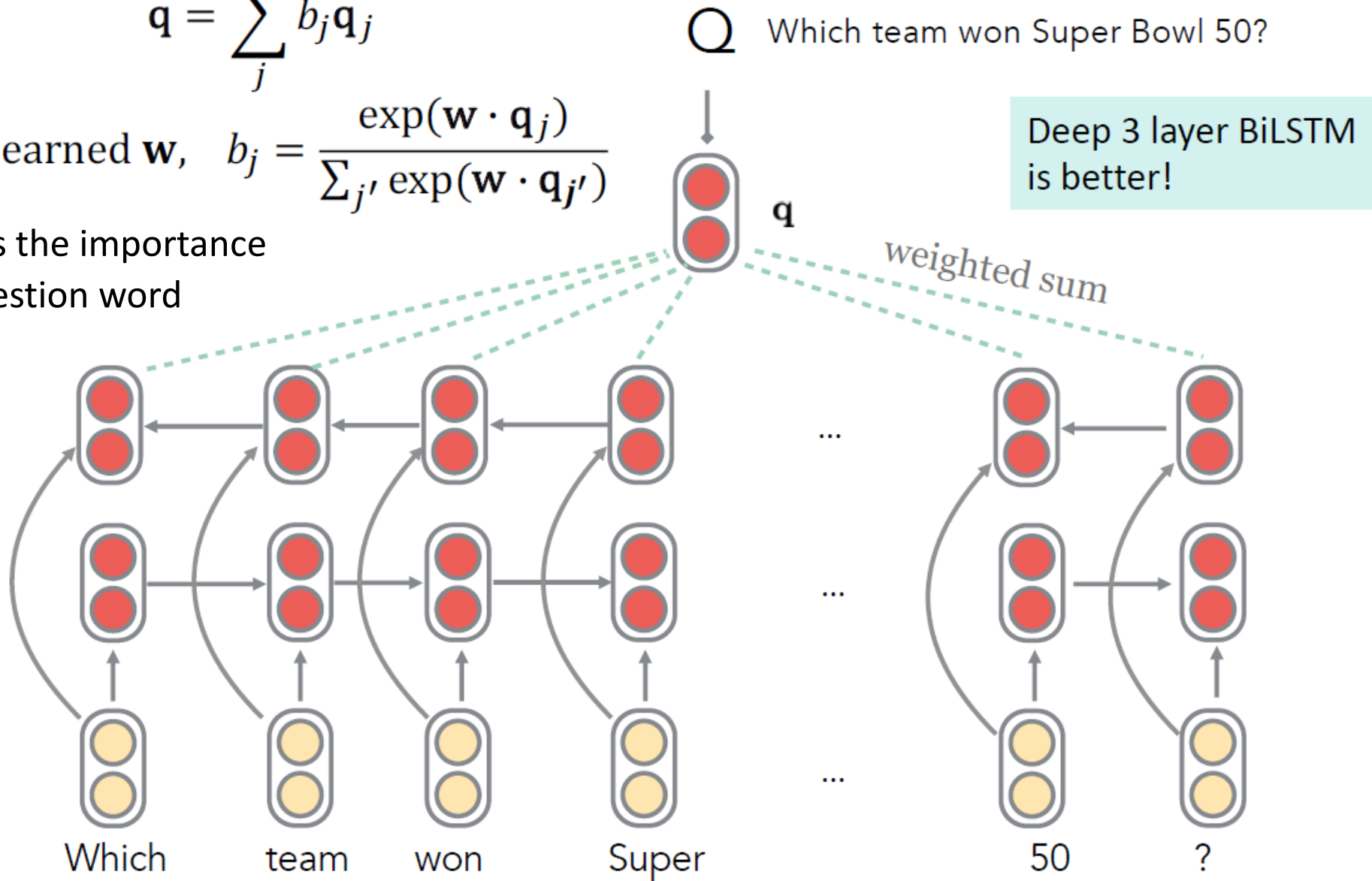
$$\mathcal{L} = - \sum \log P^{(start)}(a_{start}) - \sum \log P^{(end)}(a_{end})$$

Stanford Attentive Reader++ (Chen et al., 2016; Chen et al., 2017)

$$\mathbf{q} = \sum_j b_j \mathbf{q}_j$$

For learned \mathbf{w} ,
$$b_j = \frac{\exp(\mathbf{w} \cdot \mathbf{q}_j)}{\sum_{j'} \exp(\mathbf{w} \cdot \mathbf{q}_{j'})}$$

b_j : Encodes the importance of each question word



Stanford Attentive Reader++

- p_i : Vector representation of each token in passage

Made from concatenation of Word embedding (GloVe 300d) and

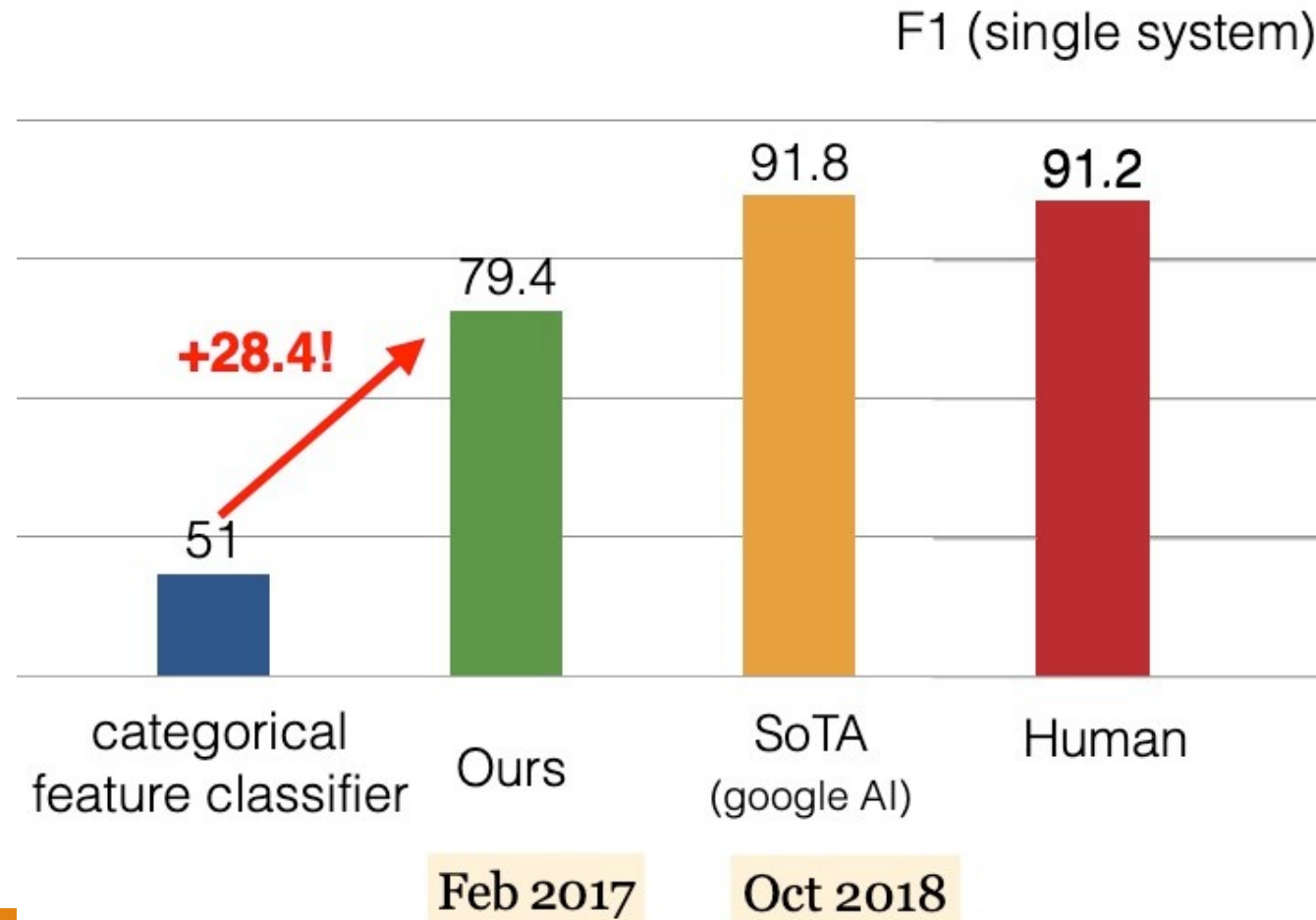
- Linguistic features: POS & NER tags, one-hot encoded
- Term frequency (unigram probability)
- Exact match: whether the word appears in the question
 - 3 binary features: exact, uncased, lemma
- Aligned question embedding: 'soft' similarity (e.g. "car" vs "vehicle")

$$f_{align}(p_i) = \sum_j a_{i,j} E(q_j) \quad a_{i,j} = \frac{\exp(\alpha(E(p_i)) \cdot \alpha(E(q_j)))}{\sum_{j'} \exp(\alpha(E(p_i)) \cdot \alpha(E(q_{j'})))}$$

Where α is a simple one layer FFNN

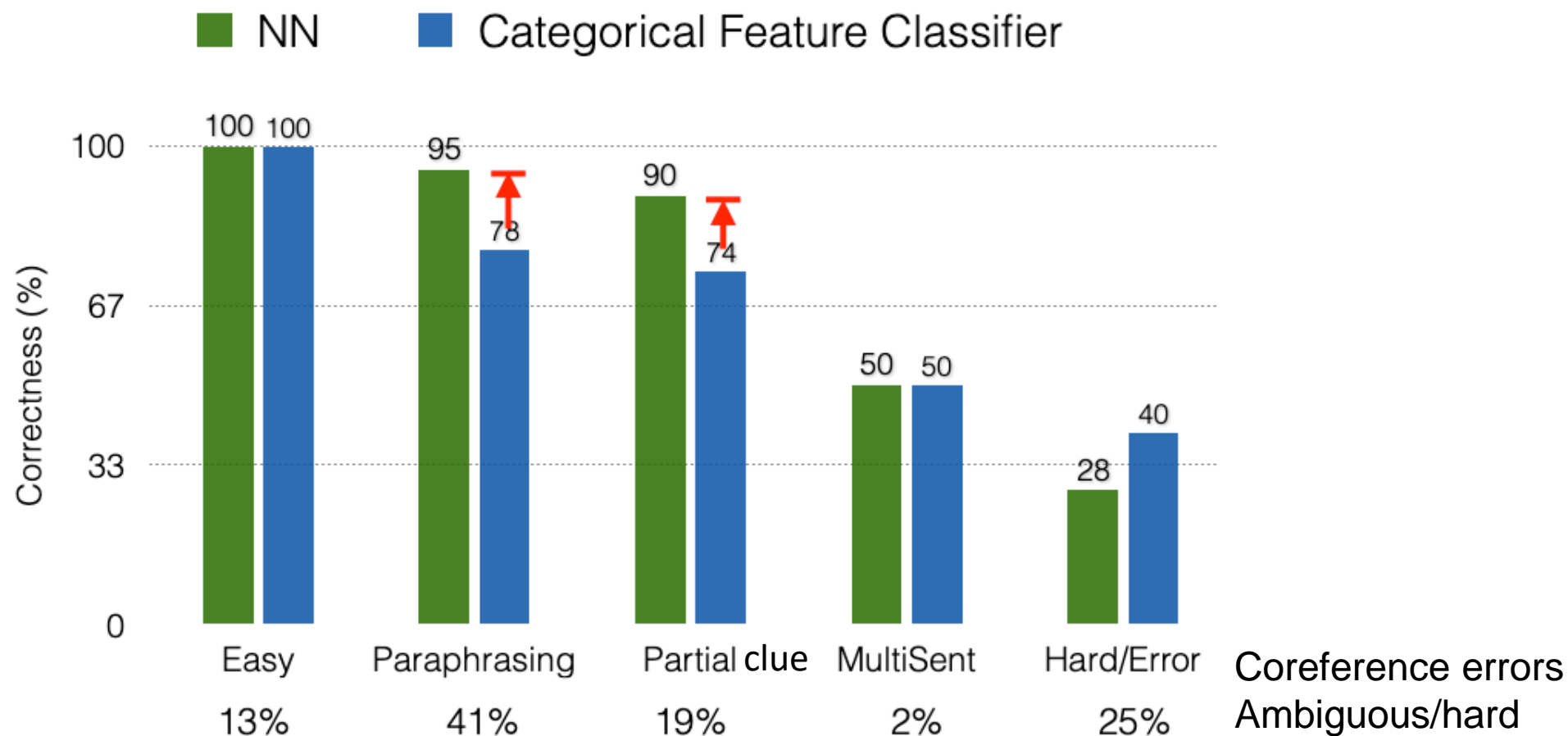
E: Glove embedding

A big win for neural models



What do these neural models do?

(Chen et al., 2016)



DrQA Demo

Hi!



Hello! Please ask a question.

What is question answering?



a computer science discipline within the fields of information retrieval and natural language processing

Who was the winning pitcher in the 1956 World Series?



Don Larsen

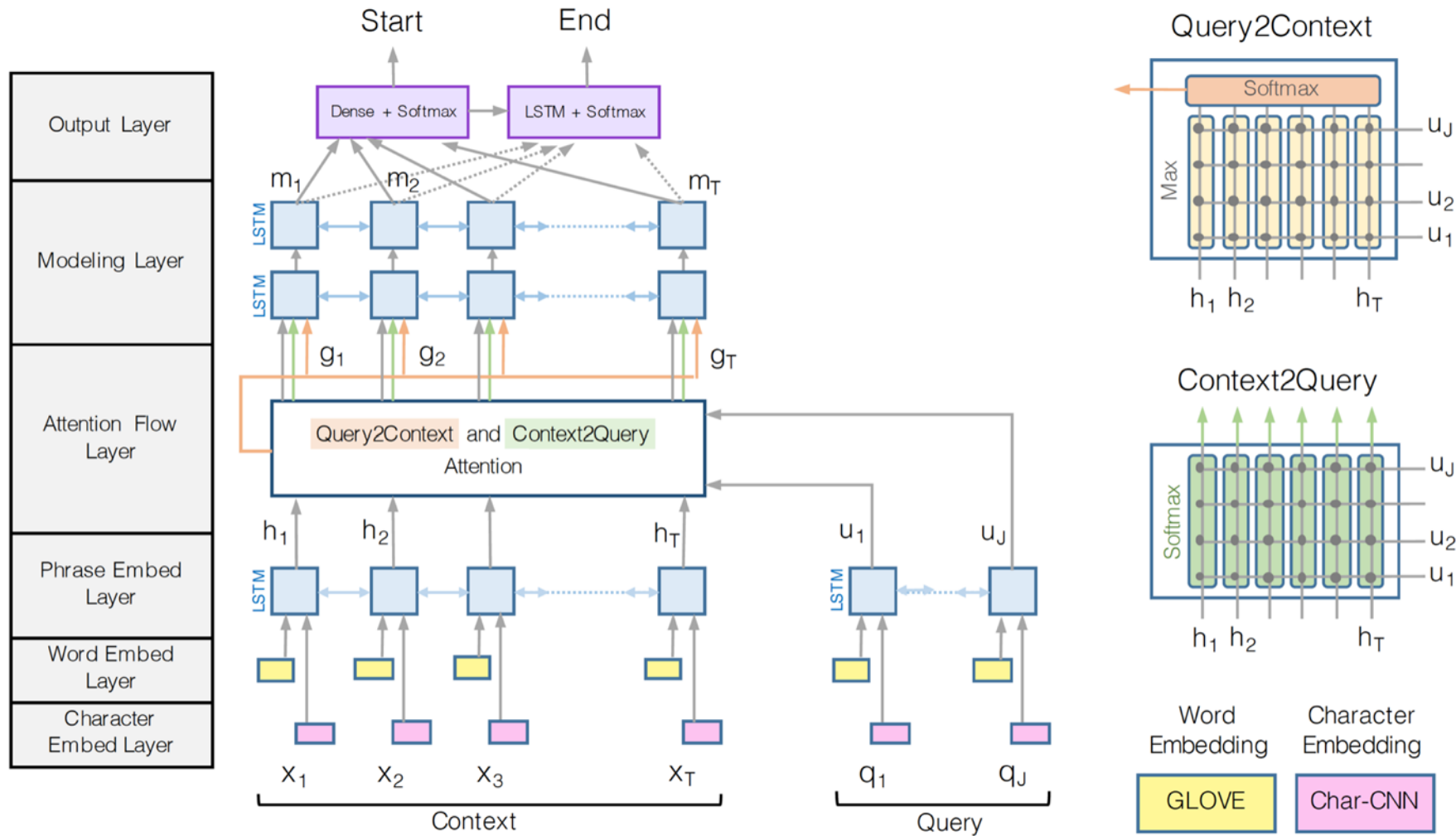
What is the answer to life, the universe, and everything?



42

4. BiDAF: Bi-Directional Attention Flow for Machine Comprehension

(Seo, Kembhavi, Farhadi, Hajishirzi, ICLR 2017)



BiDAF

- There are variants of and improvements to the BiDAF architecture over the years, but the central idea is the **Attention Flow layer**
- Idea: attention should flow both ways – from the context to the question and from the question to the context

- **Make similarity matrix** (with \mathbf{w} of dimension 6d): $S_{ij} = \mathbf{w}_{\text{sim}}^T [\mathbf{c}_i; \mathbf{q}_j; \mathbf{c}_i \circ \mathbf{q}_j] \in \mathbb{R}$

- Context-to-Question (C2Q) attention:

- (which query words are most relevant to each context word)

$$\alpha^i = \text{softmax}(\mathbf{S}_{i,:}) \in \mathbb{R}^M \quad \forall i \in \{1, \dots, N\}$$

$$\mathbf{a}_i = \sum_{j=1}^M \alpha_j^i \mathbf{q}_j \in \mathbb{R}^{2h} \quad \forall i \in \{1, \dots, N\}$$

Use this to do attention

◦ is elementwise multiplication



BiDAF

- **Attention Flow Idea:** attention should flow both ways – from the context to the question and from the question to the context
- Question-to-Context (Q2C) attention:
 - (the weighted sum of the most important words in the context with respect to the query – slight asymmetry through max)

Single representation
of the whole
context

$$\mathbf{m}_i = \max_j \mathbf{S}_{ij} \in \mathbb{R} \quad \forall i \in \{1, \dots, N\}$$

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

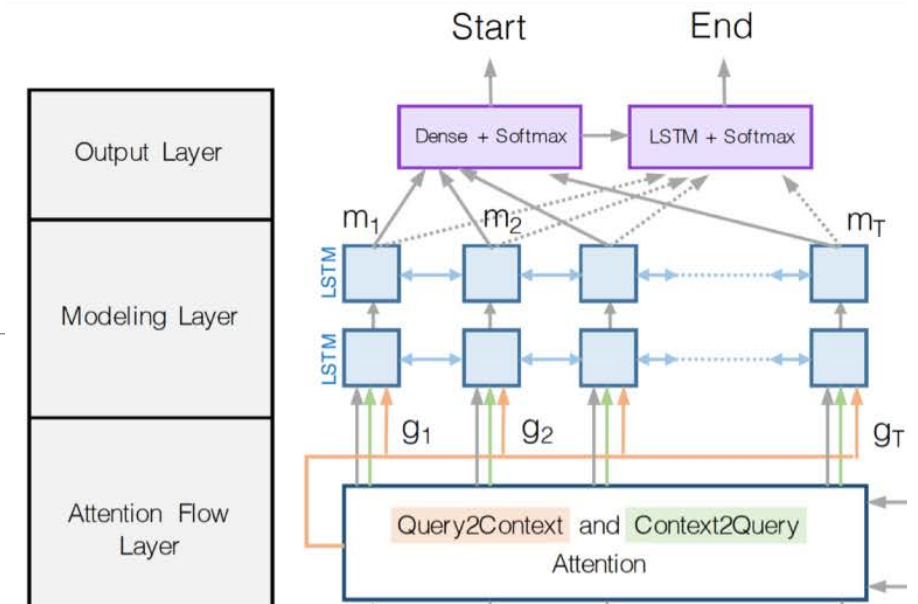
$$\mathbf{c}' = \sum_{i=1}^N \beta_i \mathbf{c}_i \in \mathbb{R}^{2h}$$

- For each passage position, output of BiDAF layer is:

$$\mathbf{b}_i = [\mathbf{c}_i; \mathbf{a}_i; \mathbf{c}_i \circ \mathbf{a}_i; \mathbf{c}_i \circ \mathbf{c}'] \in \mathbb{R}^{8h} \quad \forall i \in \{1, \dots, N\}$$

BiDAF

- There is then a “modelling” layer:
 - Another deep (2-layer) BiLSTM over the passage
- And answer span selection is more complex:
 - Start: Pass outputs of BiDAF attention flow layer (G) and of modelling layer (M), concatenated, to a **dense FF layer and then a softmax**
 - End: Put output of modelling layer (M) through another BiLSTM to give M_2 and then concatenate with BiDAF attention flow layer (G) and again put through dense FF layer and a softmax



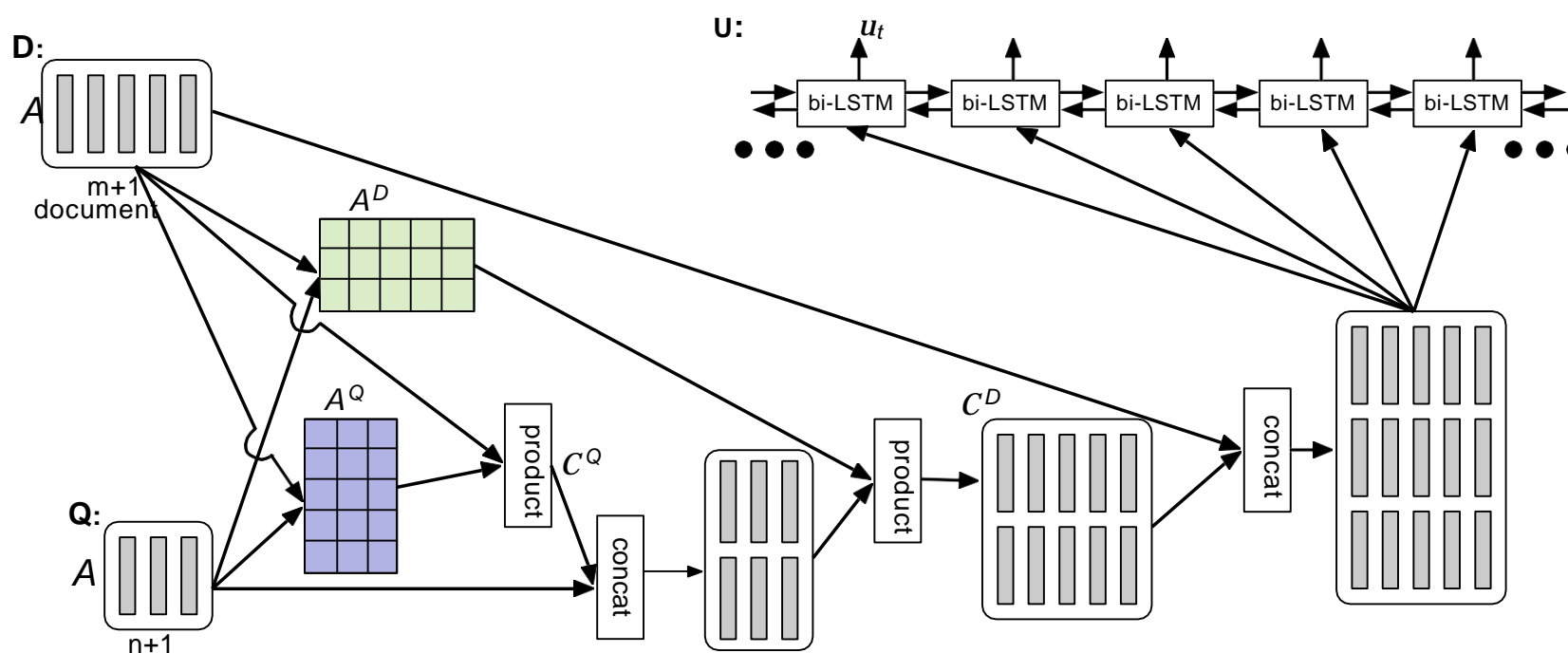
5. Recent, more advanced architectures

- Most of the work in 2016, 2017, and 2018 employed progressively more complex architectures with a multitude of variants of attention – often yielding good task gains

Dynamic Coattention Networks for Question Answering

(Caiming Xiong, Victor Zhong, ICLR 2017)

(Caiming Xiong, Victor Zhong, Richard Socher
ICLR 2017)



A^Q : question2document $C^Q = DA^Q$
 A^D : document2question

Coattention layer

- Coattention layer again provides a two-way attention between the context and the question
- However, coattention involves a second-level attention computation:
 - attending over representations that are themselves attention outputs

Co-attention: Results on SQUAD Competition

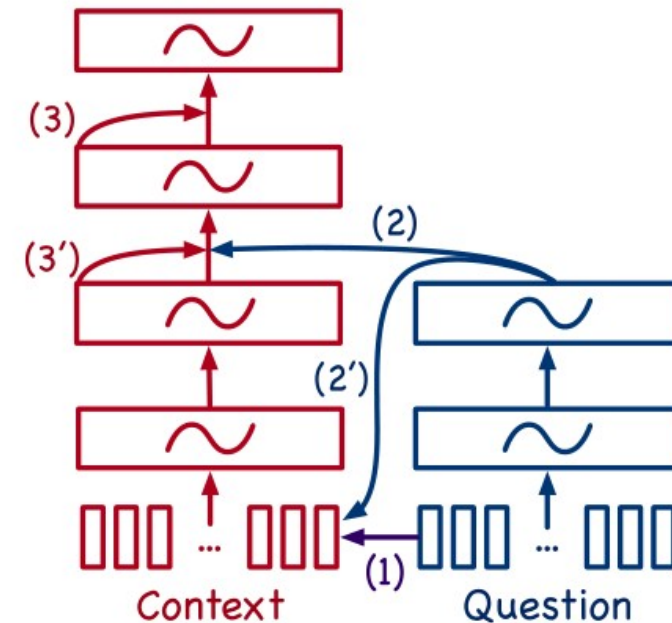
Model	Dev EM	Dev F1	Test EM	Test F1
<i>Ensemble</i>				
DCN (Ours)	70.3	79.4	71.2	80.4
Microsoft Research Asia *	—	—	69.4	78.3
Allen Institute *	69.2	77.8	69.9	78.1
Singapore Management University *	67.6	76.8	67.9	77.0
Google NYC *	68.2	76.7	—	—
<i>Single model</i>				
DCN (Ours)	65.4	75.6	66.2	75.9
Microsoft Research Asia *	65.9	75.2	65.5	75.0
Google NYC *	66.4	74.9	—	—
Singapore Management University *	—	—	64.7	73.7
Carnegie Mellon University *	—	—	62.5	73.3
Dynamic Chunk Reader (Yu et al., 2016)	62.5	71.2	62.5	71.0
Match-LSTM (Wang & Jiang, 2016)	59.1	70.0	59.5	70.3
Baseline (Rajpurkar et al., 2016)	40.0	51.0	40.4	51.0
Human (Rajpurkar et al., 2016)	81.4	91.0	82.3	91.2

Results are at time of ICLR submission. See <https://rajpurkar.github.io/SQuAD-explorer/> for latest results

Recent, more advanced architectures

- Most of the work in 2016, 2017, and 2018 employed progressively more complex architectures with a multitude of variants of attention – often yielding good task gains

Architectures	(1)	(2)	(2')	(3)	(3')
Match-LSTM (Wang and Jiang, 2016)		✓			
DCN (Xiong et al., 2017)		✓			✓
FastQA (Weissenborn et al., 2017)	✓				
FastQAExt (Weissenborn et al., 2017)	✓	✓		✓	
BiDAF (Seo et al., 2017)		✓			✓
RaSoR (Lee et al., 2016)	✓		✓		
DrQA (Chen et al., 2017)	✓				
MPCM (Wang et al., 2016)	✓	✓			
Mnemonic Reader (Hu et al., 2017)	✓	✓		✓	
R-net (Wang et al., 2017b)		✓		✓	



FusionNet combines many forms of attention

(Huang, Zhu, Shen, Chen 2017)

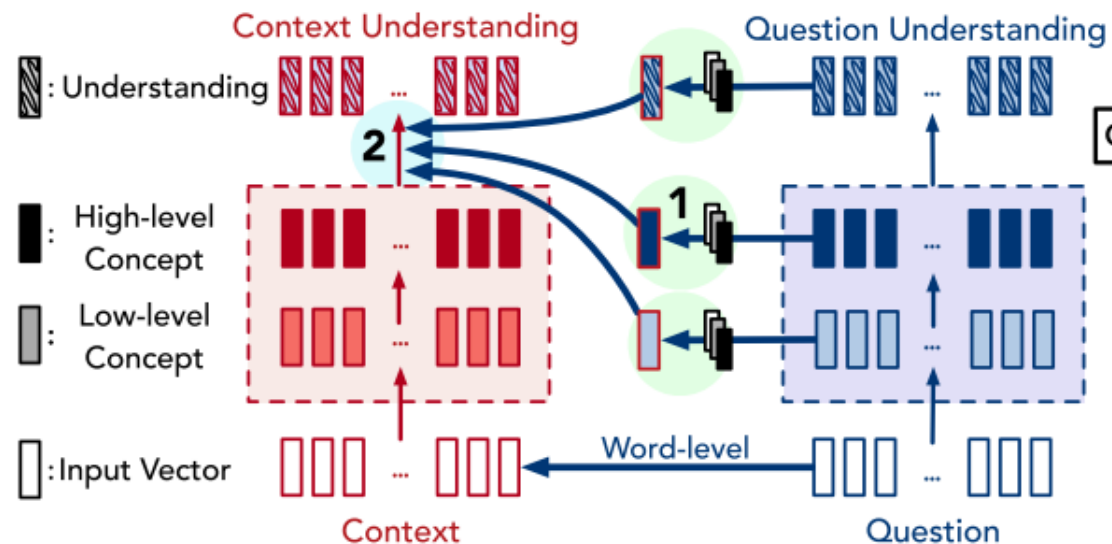
Fully-Aware Fusion Network

Fully-Aware Self-Boosted Fusion

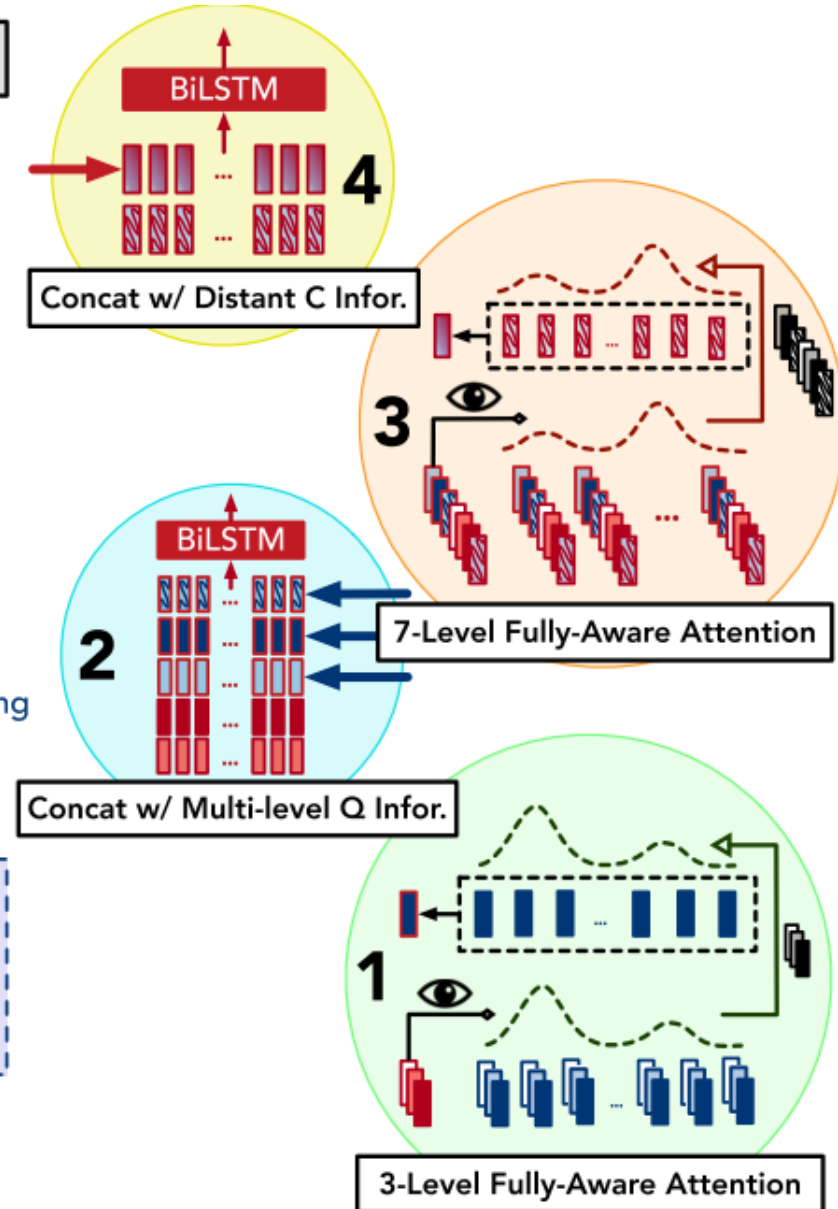


The above can be used to capture long range info.

Fully-Aware Multi-level Fusion



Each upward arrow represents one layer of BiLSTM



between context and question

7. ELMo and BERT preview

Contextual word representations

Using language model-like objectives



Elmo

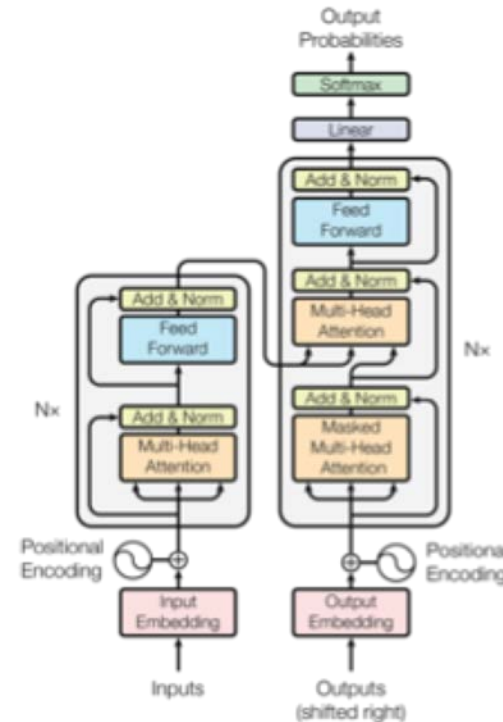
(Peters et al, 2018)



Bert

(Devlin et al, 2018)

The transformer architecture used in BERT is sort of attention on steroids. More later!



(Vaswani et al, 2017)

DrQA: Open-domain Question Answering

Q: How many of Warsaw's inhabitants spoke Polish in 1933?



Document Retriever

Dataset	Wiki Search	Doc. plain	Retriever +bigrams
SQuAD	62.7	76.1	77.8
CuratedTREC	81.0	85.2	86.0
WebQuestions	73.7	75.5	74.4
WikiMovies	61.7	54.4	70.3

Traditional
tf.idf
inverted
index +
efficient
bigram
hash

AI6122: Text Data
Management
and Processing

For **70–86%** of questions, the answer segment appears in the top 5 articles