

Lecture 14: Question Answering

Wei Xu

(many slides from Greg Durrett)

QA is very broad

- ▶ Factoid QA: *what states border Mississippi?, when was Barack Obama born?*
 - ▶ Lots of this could be handled by QA from a knowledge base, if we had a big enough knowledge base
- ▶ “Question answering” as a term is so broad as to be meaningless
 - ▶ *Is $P=NP$?*
 - ▶ *What is $4+5$?*
 - ▶ *What is the translation of [sentence] into French?* [McCann et al., 2018]

Classical Question Answering

- ▶ Form semantic representation from semantic parsing, execute against structured knowledge base

Q: “where was Barack Obama born”

$\lambda x. \text{type}(x, \text{Location}) \wedge \text{born_in}(\text{Barack_Obama}, x)$

(other representations like SQL possible too...)

- ▶ How to deal with open-domain data/relations? Need data to learn how to ground every predicate or need to be able to produce predicates in a zero-shot way

Reading Comprehension

- ▶ “AI challenge problem”: answer question given context
- ▶ Recognizing Textual Entailment (2006)
- ▶ MCTest (2013): 500 passages, 4 questions per passage
- ▶ Two questions per passage explicitly require cross-sentence reasoning

One day, James thought he would go into town and see what kind of trouble he could get into. He went to the grocery store and pulled all the pudding off the shelves and ate two jars. Then he walked to the fast food restaurant and ordered 15 bags of fries. He didn't pay, and instead headed home.

3) Where did James go after he went to the grocery store?

A) his deck

B) his freezer

C) a fast food restaurant

D) his room

Dataset Explosion

- ▶ 10+ QA datasets released since 2015
 - ▶ Children's Book Test, CNN/Daily Mail, SQuAD, TriviaQA are most well-known (others: SearchQA, MS Marco, RACE, WikiHop, ...)
- ▶ Question answering: questions are in natural language
 - ▶ Answers: multiple choice or require picking from the passage
 - ▶ Require human annotation
- ▶ "Cloze" task: word (often an entity) is removed from a sentence
 - ▶ Answers: multiple choice, pick from passage, or pick from vocabulary
 - ▶ Can be created automatically from things that aren't questions

Children's Book Test

"Well, Miss Maxwell, I think it only fair to tell you that you may have trouble with those boys when they do come. Forewarned is forearmed, you know. Mr. Cropper was opposed to our hiring you. Not, of course, that he had any personal objection to you, but he is set against female teachers, and when a Cropper is set there is nothing on earth can change him. He says female teachers can't keep order. He 's started in with a spite at you on general principles, and the boys know it. They know he'll back them up in secret, no matter what they do, just to prove his opinions. Cropper is sly and slippery, and

S: 1 Mr. Cropper was opposed to our hiring you .
2 Not , of course , that he had any personal objection to you , but he is set against female teachers , and when a Cropper is set there is nothing on earth can change him .
3 He says female teachers ca n't keep order .
4 He 's started in with a spite at you on general principles , and the boys know it .
5 They know he 'll back them up in secret , no matter what they do , just to prove his opinions .
6 Cropper is sly and slippery , and it is hard to corner him . ''
7 `` Are the boys big ? ''

Mr. Baxter privately had no hope that they would, but Esther hoped for the best. She could not believe that Mr. Cropper would carry his prejudices into a personal application. This conviction was strengthened when he overtook her walking from school the next day and drove her home. He was a big, handsome man with a very suave, polite manner. He asked interestedly about her school and her work, hoped she was getting on well, and said he had two young rascals of his own to send soon. Esther felt relieved. She thought that **????** had exaggerated matters a little.

r their age .
he trouble .
you around their fingers .
'm afraid .
ght after all . ''
that they would , but Esther hoped for the
cropper would carry his prejudices into a
when he overtook her walking from school the
a very suave , polite manner .
school and her work , hoped she was getting on
scals of his own to send soon .
exaggerated matters a little .
ngers, manner, objection, opinion, right, spite.

- ▶ Children's Book Test: take a section of a children's story, block out an entity and predict it (one-doc multi-sentence cloze task)

Hill et al. (2015)

bAbI

- ▶ Evaluation on 20 tasks proposed as building blocks for building “AI-complete” systems
- ▶ Various levels of difficulty, exhibit different linguistic phenomena
- ▶ Small vocabulary, language isn’t truly “natural”

Task 1: Single Supporting Fact

Mary went to the bathroom.
John moved to the hallway.
Mary travelled to the office.
Where is Mary? **A:office**

Task 2: Two Supporting Facts

John is in the playground.
John picked up the football.
Bob went to the kitchen.
Where is the football? **A:playground**

Task 13: Compound Coreference

Daniel and Sandra journeyed to the office.
Then they went to the garden.
Sandra and John travelled to the kitchen.
After that they moved to the hallway.
Where is Daniel? **A: garden**

Task 14: Time Reasoning

In the afternoon Julie went to the park.
Yesterday Julie was at school.
Julie went to the cinema this evening.
Where did Julie go after the park? **A:cinema**
Where was Julie before the park? **A:school**

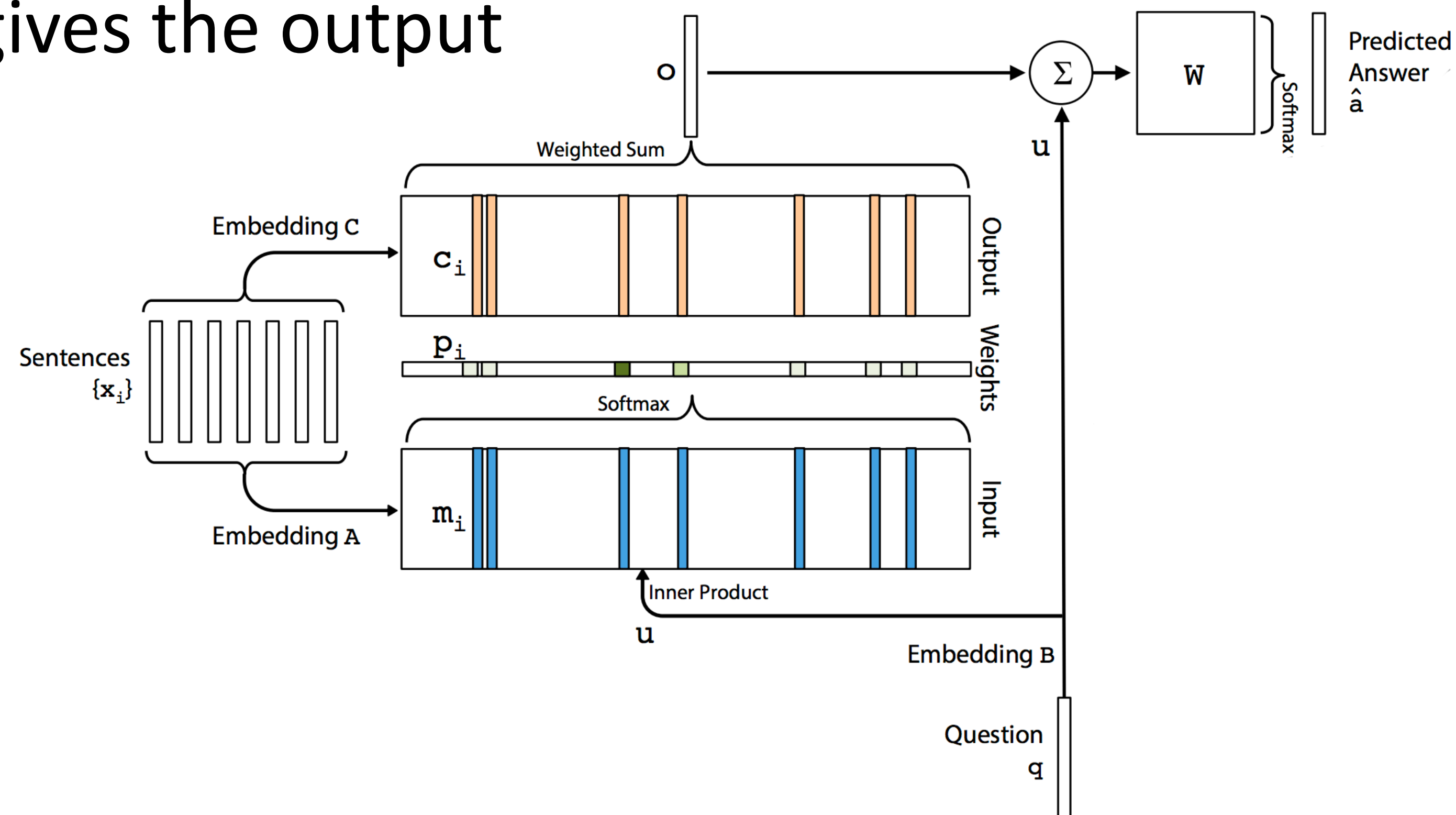
Dataset Properties

- ▶ Axis 1: QA vs. cloze (Children's Book Test)
- ▶ Axis 2: single-sentence vs. passage
 - ▶ Often shallow methods work well because most answers are in a single sentence (SQuAD, MCTest)
 - ▶ Some explicitly require linking between multiple sentences (MCTest)
- ▶ Axis 3: single-document (datasets in this lecture) vs. multi-document (TriviaQA, WikiHop, HotPotQA, ...)

Memory Networks

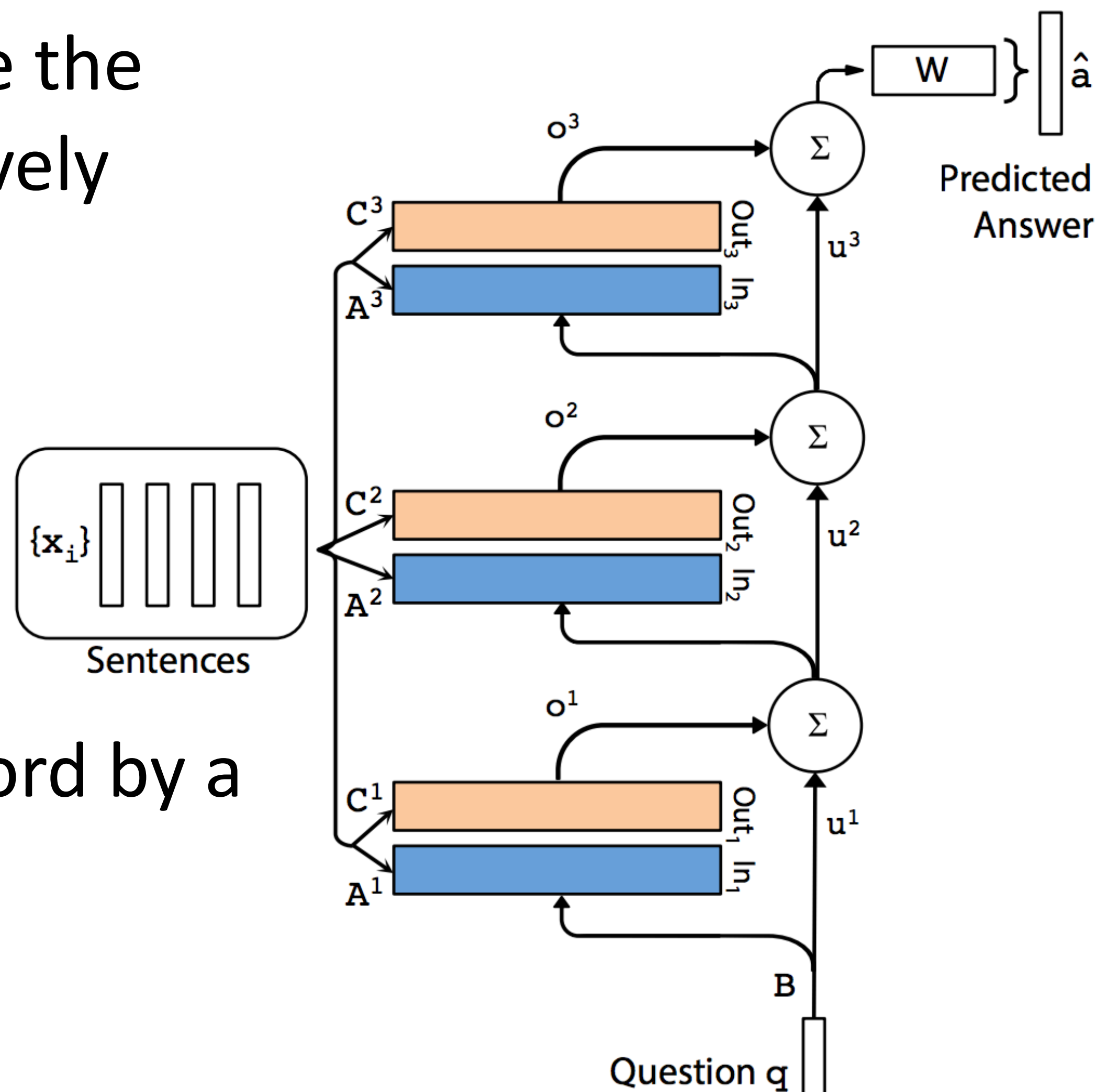
Memory Networks

- ▶ Memory networks let you reference input with attention
- ▶ Encode input items into two vectors: a **key** and a **value**
- ▶ Keys compute attention weights given a query, weighted sum of values gives the output



Memory Networks

- ▶ Three layers of memory network where the query representation is updated additively based on the memories at each step
- ▶ How to encode the sentences?
 - ▶ Bag of words (average embeddings)
 - ▶ Positional encoding: multiply each word by a vector capturing position in sentence



(b)

Evaluation: bAbI

Task	Baseline				MemN2N			
	Strongly Supervised MemNN [22]	LSTM [22]	MemNN WSH	BoW	PE	1 hop PE LS joint	2 hops PE LS joint	3 hops PE LS joint
Mean error (%)	6.7	51.3	40.2	25.1	20.3	25.8	15.6	13.3
Failed tasks (err. > 5%)	4	20	18	15	13	17	11	11

- 3-hop memory network does pretty well, better than LSTM at processing these types of examples

Story (16: basic induction)	Support	Hop 1	Hop 2	Hop 3
Brian is a frog.	yes	0.00	0.98	0.00
Lily is gray.		0.07	0.00	0.00
Brian is yellow.	yes	0.07	0.00	1.00
Julius is green.		0.06	0.00	0.00
Greg is a frog.	yes	0.76	0.02	0.00
What color is Greg? Answer: yellow Prediction: yellow				

Evaluation: Children's Book Test

METHODS	NAMED ENTITIES
HUMANS (QUERY) ^(*)	0.520
HUMANS (CONTEXT+QUERY) ^(*)	0.816
MAXIMUM FREQUENCY (CORPUS)	0.120
MAXIMUM FREQUENCY (CONTEXT)	0.335
SLIDING WINDOW	0.168
WORD DISTANCE MODEL	0.398
KNESER-NEY LANGUAGE MODEL	0.390
KNESER-NEY LANGUAGE MODEL + CACHE	0.439
LSTMs (QUERY)	0.408
LSTMs (CONTEXT+QUERY)	0.418
CONTEXTUAL LSTMs (WINDOW CONTEXT)	0.436
MEMNNS (LEXICAL MEMORY)	0.431
MEMNNS (WINDOW MEMORY)	0.493
MEMNNS (SENTENTIAL MEMORY + PE)	0.318
MEMNNS (WINDOW MEMORY + SELF-SUP.)	0.666

▶ Outperforms LSTMs substantially with the right supervision

Memory Network Takeaways

- ▶ Memory networks provide a way of attending to abstractions over the input
- ▶ Useful for cloze tasks where far-back context is necessary
- ▶ What can we do with more basic attention?

CNN/Daily Mail: Attentive Reader

CNN/Daily Mail

- ▶ Single-document, (usually) single-sentence cloze task
- ▶ Formed based on article summaries — information should mostly be present, makes it easier than Children's Book Test
- ▶ Need to process the question, can't just use LSTM LMs

Passage

(@entity4) if you feel a ripple in the force today , it may be the news that the official @entity6 is getting its first gay character . according to the sci-fi website @entity9 , the upcoming novel " @entity11 " will feature a capable but flawed @entity13 official named @entity14 who " also happens to be a lesbian . " the character is the first gay figure in the official @entity6 -- the movies , television shows , comics and books approved by @entity6 franchise owner @entity22 -- according to @entity24 , editor of " @entity6 " books at @entity28 imprint @entity26 .

Question

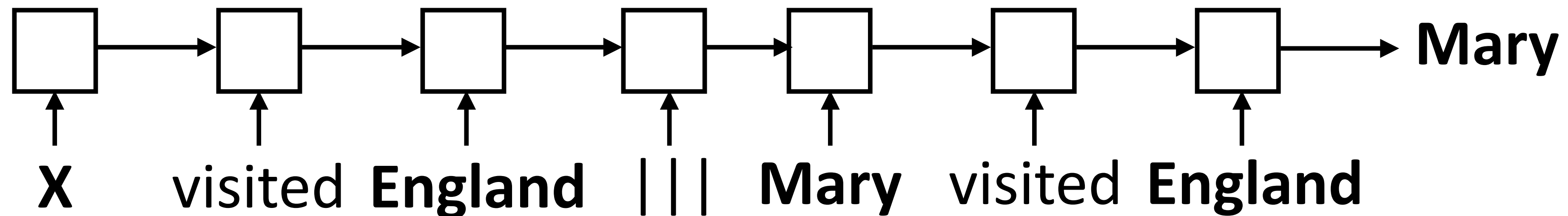
characters in " @placeholder " movies have gradually become more diverse

Answer

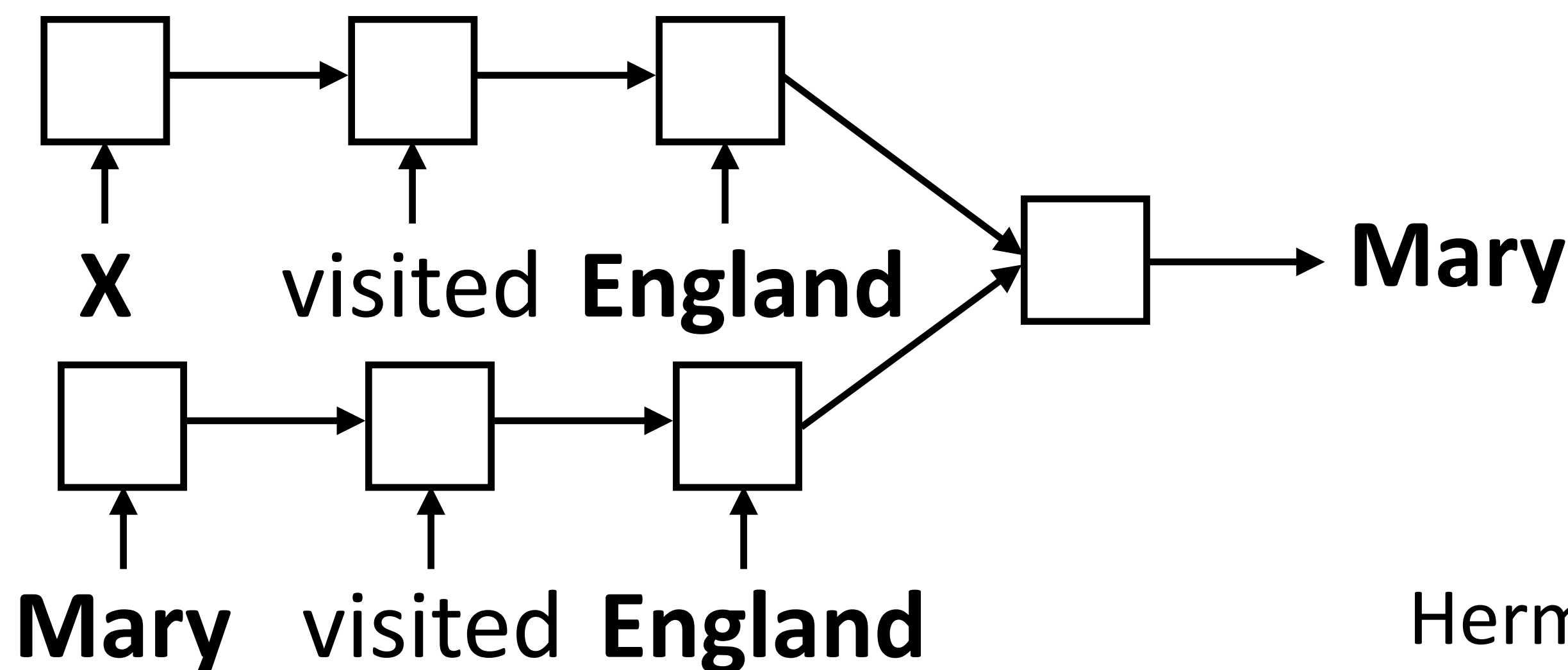
@entity6

CNN/Daily Mail

- ▶ LSTM reader: encode question, encode passage, predict entity



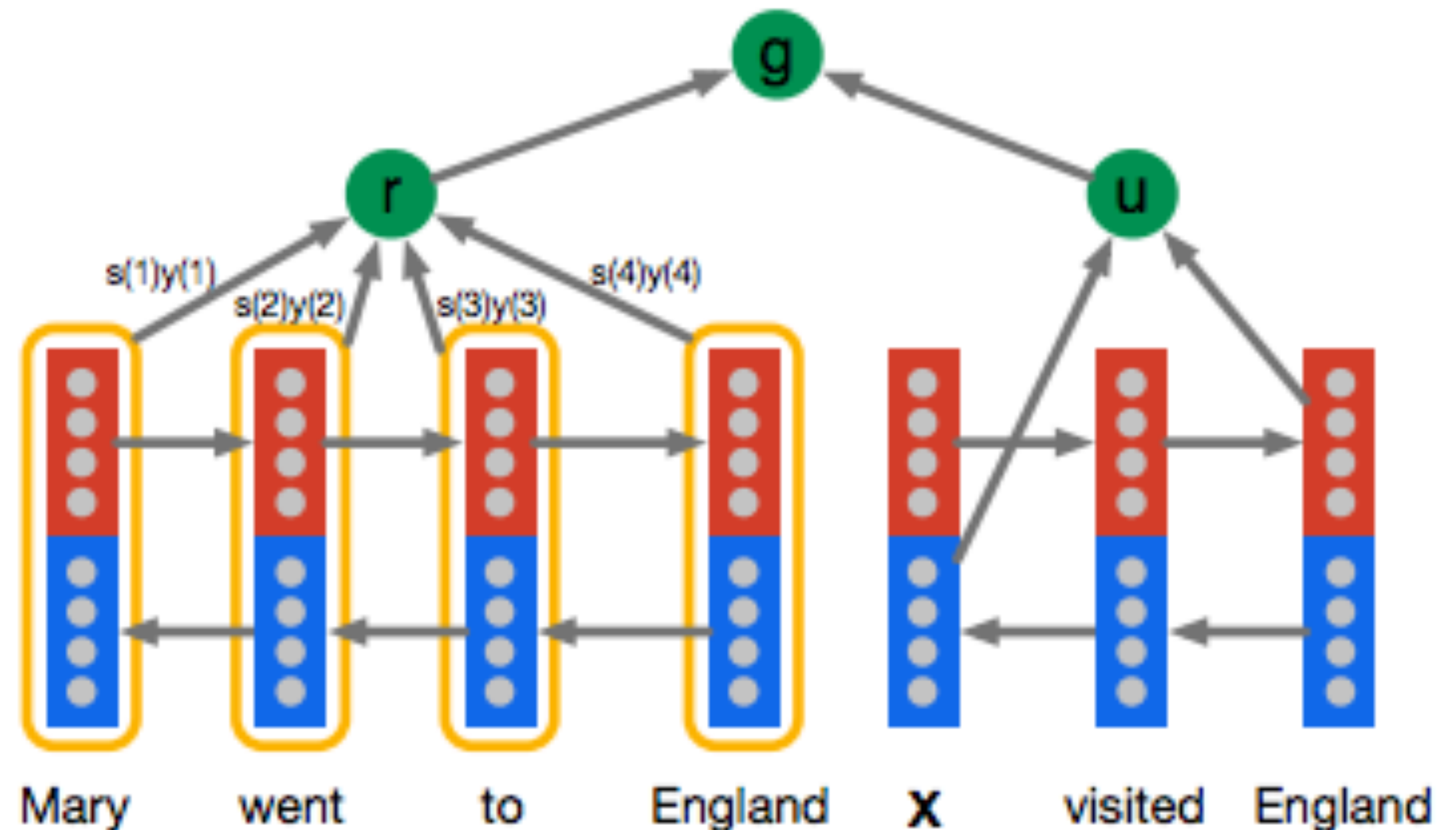
- ▶ Can also use textual entailment-like models



Multiclass classification problem over entities in the document

CNN/Daily Mail

- ▶ Attentive reader:
 - u = encode query
 - s = encode sentence
 - $r = \text{attention}(u \rightarrow s)$
 - $\text{prediction} = f(\text{candidate}, u, r)$
- ▶ Uses fixed-size representations for the final prediction, multiclass classification



CNN/Daily Mail

- ▶ Chen et al (2016): small changes to the attentive reader
- ▶ Additional analysis of the task found that many of the remaining questions were unanswerable or extremely difficult

	CNN		Daily Mail	
	valid	test	valid	test
Maximum frequency	30.5	33.2	25.6	25.5
Exclusive frequency	36.6	39.3	32.7	32.8
Frame-semantic model	36.3	40.2	35.5	35.5
Word distance model	50.5	50.9	56.4	55.5
Deep LSTM Reader	55.0	57.0	63.3	62.2
Uniform Reader	39.0	39.4	34.6	34.4
Attentive Reader	61.6	63.0	70.5	69.0
Stanford Attentive Reader	76.2	76.5	79.5	78.7

SQuAD: Bidirectional Attention Flow

SQuAD

- ▶ Single-document, single-sentence question-answering task where the answer is always a substring of the passage
- ▶ Predict start and end indices of the answer in the passage

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.

Question: Which NFL team won Super Bowl 50?

Answer: Denver Broncos

Question: What does AFC stand for?

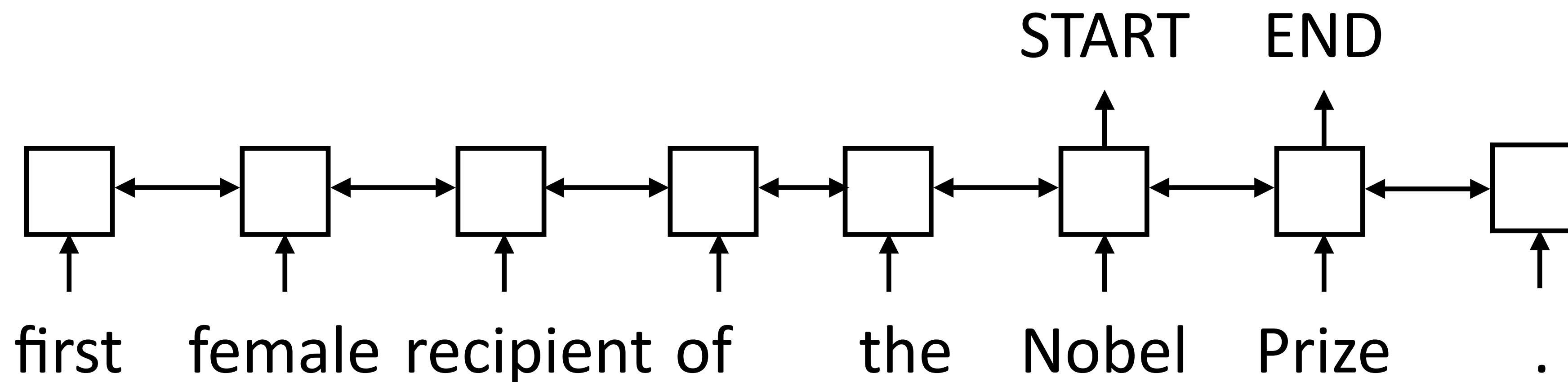
Answer: American Football Conference

Question: What year was Super Bowl 50?

Answer: 2016

SQuAD

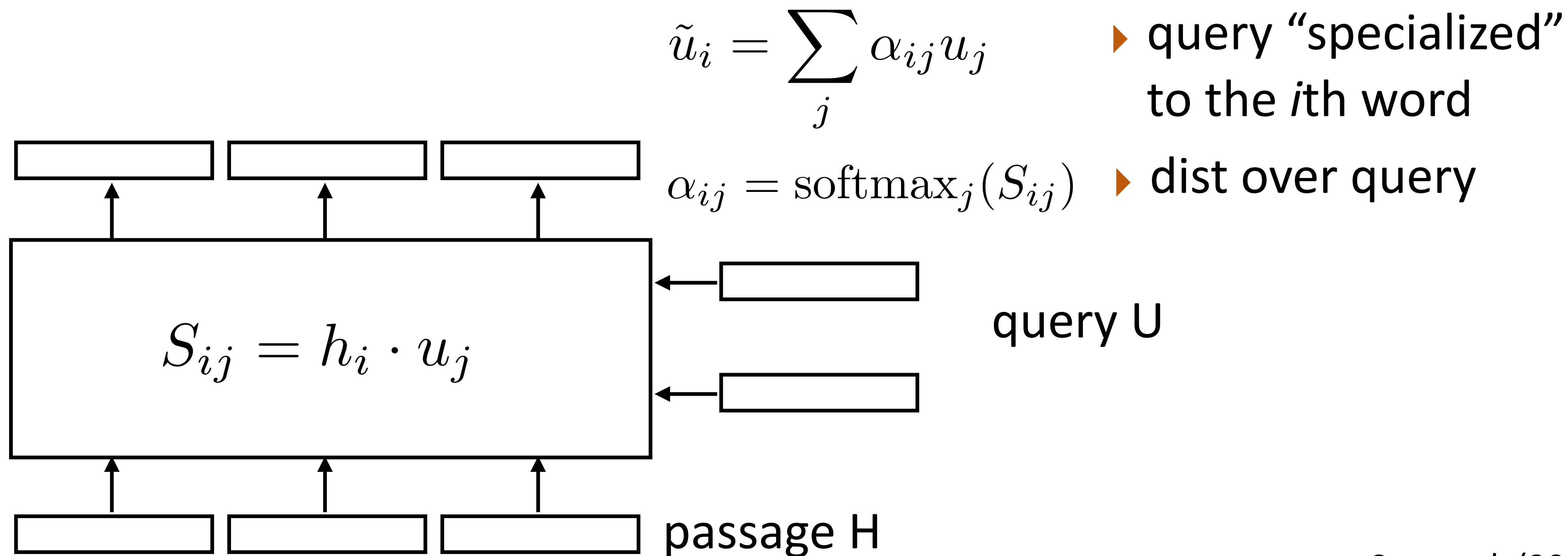
What was Marie Curie the first female recipient of?



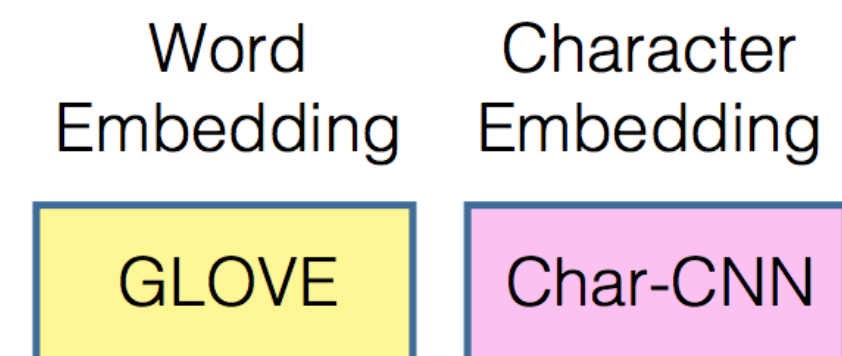
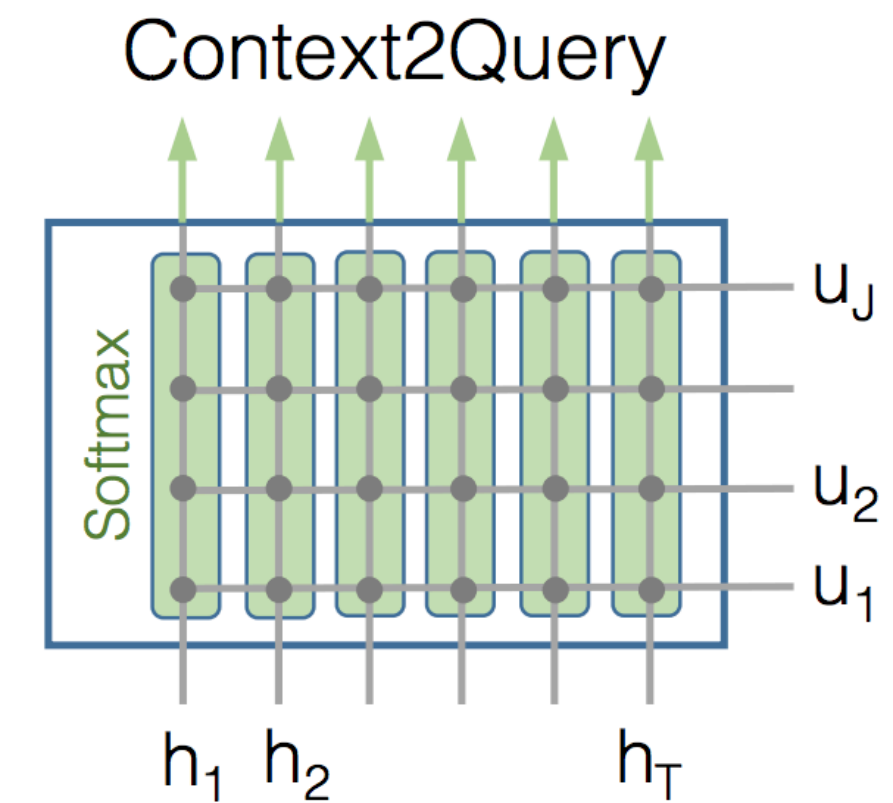
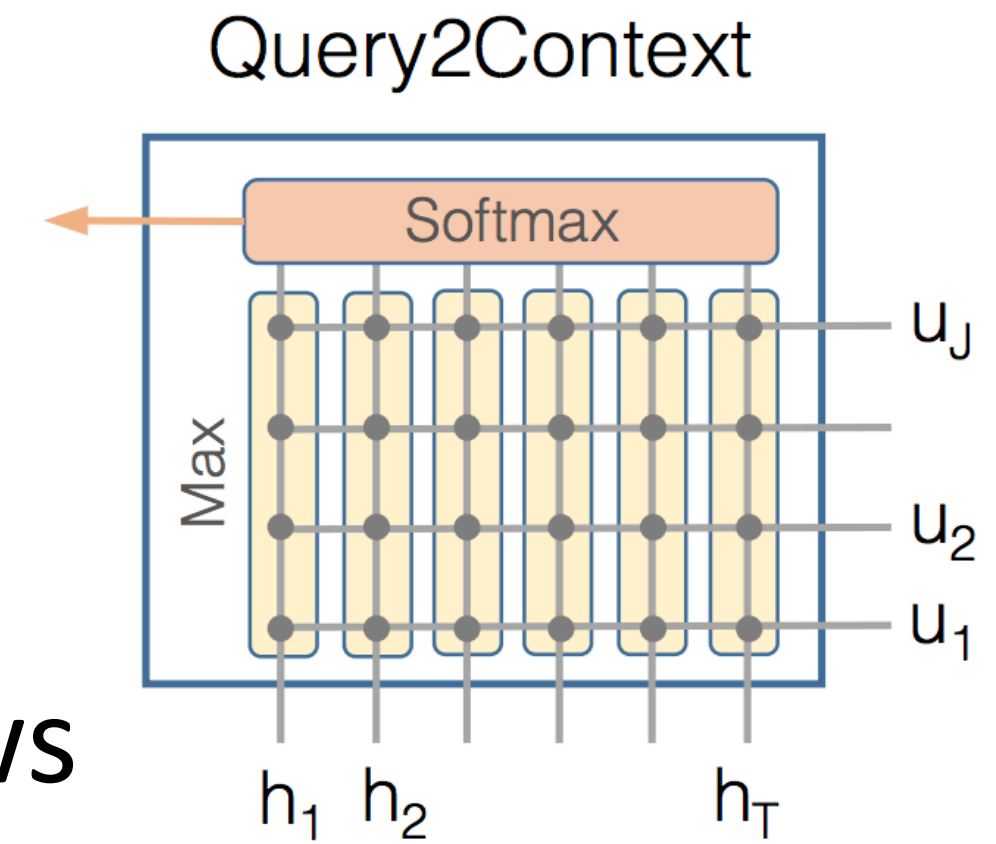
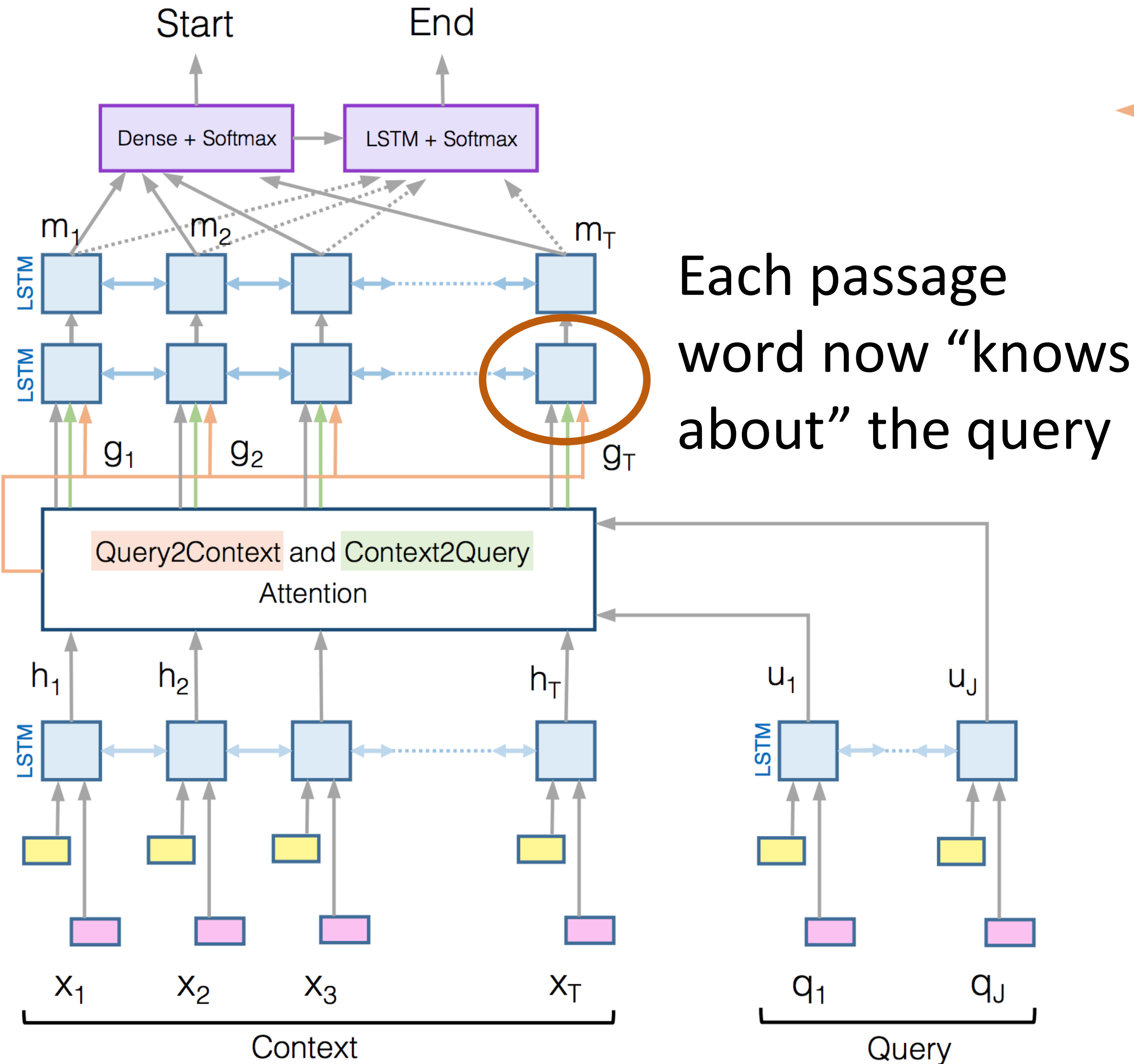
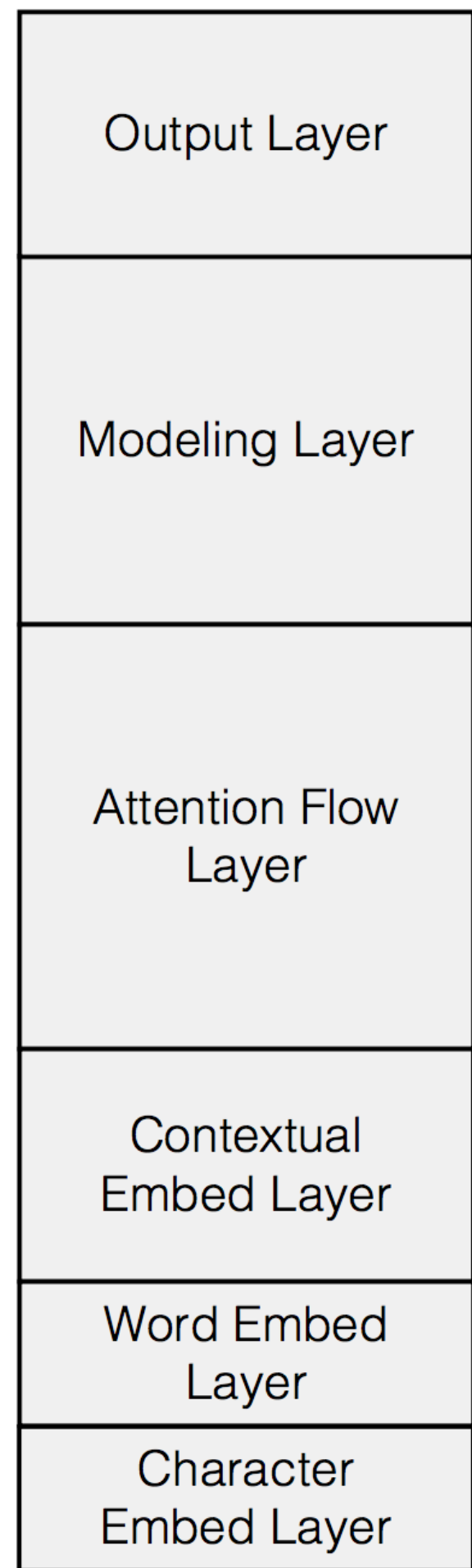
- ▶ Like a tagging problem over the sentence (not multiclass classification), but we need some way of attending to the query

Bidirectional Attention Flow

- ▶ Passage (context) and query are both encoded with BiLSTMs
- ▶ Context-to-query attention: compute softmax over columns of S , take weighted sum of u based on attention weights for each passage word

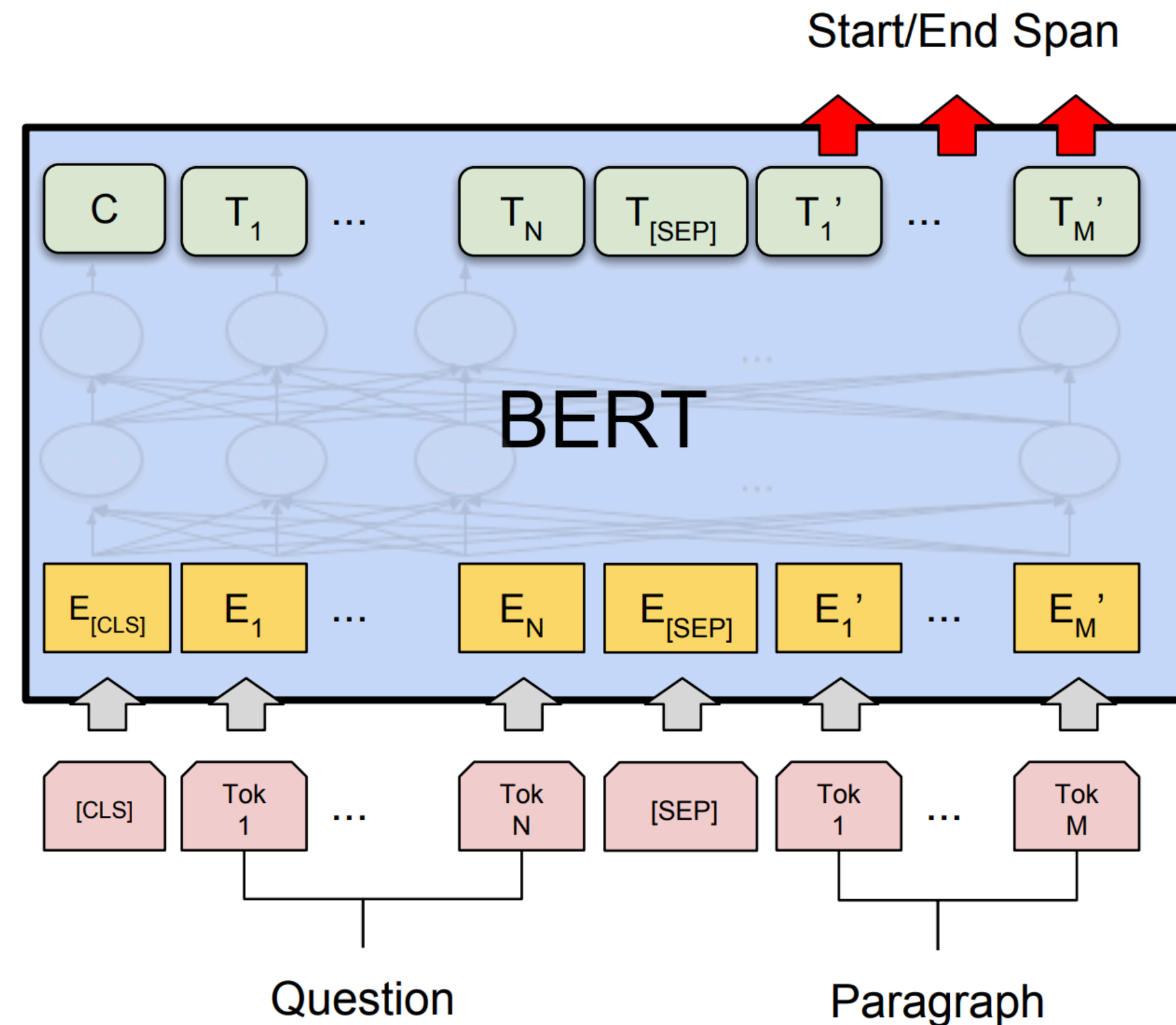


Bidirectional Attention Flow





QA with BERT



What was Marie Curie the first female recipient of ? [SEP] One of the most famous people born in Warsaw was Marie ...

- ▶ Predict start and end positions in passage
- ▶ No need for cross-attention mechanisms!

SQuAD SOTA: 2018

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

- ▶ BiDAF: 73 EM / 81 F1
- ▶ nlnet, QANet, r-net — dueling super complex systems (much more than BiDAF...)
- ▶ BERT: transformer-based approach with pretraining on 3B tokens

SQuAD 2.0 SOTA: Spring 2019

Rank	Model	EM	F1
	Human Performance <i>Stanford University</i> (Rajpurkar & Jia et al. '18)	86.831	89.452
1 Mar 20, 2019	BERT + DAE + AoA (ensemble) <i>Joint Laboratory of HIT and iFLYTEK Research</i>	87.147	89.474
2 Mar 15, 2019	BERT + ConvLSTM + MTL + Verifier (ensemble) <i>Layer 6 AI</i>	86.730	89.286
3 Mar 05, 2019	BERT + N-Gram Masking + Synthetic Self-Training (ensemble) <i>Google AI Language</i> https://github.com/google-research/bert	86.673	89.147
4 Apr 13, 2019	SemBERT(ensemble) <i>Shanghai Jiao Tong University</i>	86.166	88.886
5 Mar 16, 2019	BERT + DAE + AoA (single model) <i>Joint Laboratory of HIT and iFLYTEK Research</i>	85.884	88.621
6 Mar 05, 2019	BERT + N-Gram Masking + Synthetic Self-Training (single model) <i>Google AI Language</i> https://github.com/google-research/bert	85.150	87.715
7 Jan 15, 2019	BERT + MMFT + ADA (ensemble) <i>Microsoft Research Asia</i>	85.082	87.615

- ▶ SQuAD 2.0: harder dataset because some questions are unanswerable
- ▶ Industry contest

SQuAD 2.0 SOTA: Fall 2019

Rank	Model	EM	F1
	Human Performance <i>Stanford University</i> (Rajpurkar & Jia et al. '18)	86.831	89.452
1 Sep 18, 2019	ALBERT (ensemble model) <i>Google Research & TTIC</i> https://arxiv.org/abs/1909.11942	89.731	92.215
2 Jul 22, 2019	XLNet + DAAF + Verifier (ensemble) <i>PINGAN Omni-Sinitic</i>	88.592	90.859
2 Sep 16, 2019	ALBERT (single model) <i>Google Research & TTIC</i> https://arxiv.org/abs/1909.11942	88.107	90.902
2 Jul 26, 2019	UPM (ensemble) <i>Anonymous</i>	88.231	90.713
3 Aug 04, 2019	XLNet + SG-Net Verifier (ensemble) <i>Shanghai Jiao Tong University & CloudWalk</i> https://arxiv.org/abs/1908.05147	88.174	90.702
4 Aug 04, 2019	XLNet + SG-Net Verifier++ (single model) <i>Shanghai Jiao Tong University & CloudWalk</i> https://arxiv.org/abs/1908.05147	87.238	90.071

▶ Performance is very saturated

▶ Harder QA settings are needed!

SQuAD 2.0 SOTA: Today

Rank	Model	EM	F1
	Human Performance <i>Stanford University</i> (Rajpurkar & Jia et al. '18)	86.831	89.452
1 Feb 21, 2021	FPNet (ensemble) <i>Ant Service Intelligence Team</i>	90.871	93.183
2 Feb 24, 2021	IE-Net (ensemble) <i>RICOH_SRCB_DML</i>	90.758	93.044
3 Apr 06, 2020	SA-Net on Albert (ensemble) <i>QIANXIN</i>	90.724	93.011
4 May 05, 2020	SA-Net-V2 (ensemble) <i>QIANXIN</i>	90.679	92.948
4 Apr 05, 2020	Retro-Reader (ensemble) <i>Shanghai Jiao Tong University</i> http://arxiv.org/abs/2001.09694	90.578	92.978
4 Feb 05, 2021	FPNet (ensemble) <i>YuYang</i>	90.600	92.899
5 Dec 01, 2020	EntitySpanFocusV2 (ensemble) <i>RICOH_SRCB_DML</i>	90.521	92.824
5 Jul 31, 2020	ATRLP+PV (ensemble) <i>Hithink RoyalFlush</i>	90.442	92.877
5 May 04, 2020	ELECTRA+ALBERT+EntitySpanFocus (ensemble) <i>SRCB_DML</i>	90.442	92.839

▶ Performance is very saturated

▶ Harder QA settings are needed!

TriviaQA

- ▶ Totally figuring this out is very challenging
- ▶ Coref:
the failed campaign movie of the same name
- ▶ Lots of surface clues:
1961, campaign, etc.
- ▶ Systems can do well without really understanding the text

Question: The Dodecanese **Campaign** of WWII that was an attempt by the Allied forces to capture islands in the Aegean Sea was the inspiration for which acclaimed 1961 commando film?

Answer: The Guns of Navarone

Excerpt: The Dodecanese Campaign of World War II was an attempt by Allied forces to capture the Italian-held Dodecanese islands in the Aegean Sea following the surrender of Italy in September 1943, and use them as bases against the German-controlled Balkans. The **failed campaign**, and in particular the Battle of Leros, inspired the 1957 novel **The Guns of Navarone** and the successful **1961 movie of the same name.**

What are these models learning?

- ▶ “Who...”: knows to look for people
- ▶ “Which film...”: can identify movies and then spot keywords that are related to the question
- ▶ Unless questions are made super tricky (target closely-related entities who are easily confused), they’re usually not so hard to answer

Takeaways

- ▶ Many flavors of reading comprehension tasks: cloze or actual questions, single or multi-sentence
- ▶ Memory networks let you reference input in an attention-like way, useful for generalizing language models to long-range reasoning
- ▶ Complex attention schemes can match queries against input texts and identify answers