

Lecture 14: Question Answering I

Alan Ritter

(many slides from Greg Durrett)

Classical Question Answering

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

Classical Question Answering

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

Q: “where was Barack Obama born”

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...)

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

QA from Open IE

(a) CCG parse builds an underspecified semantic representation of the sentence.

Former	municipalities	in	Brandenburg
N/N $\lambda f \lambda x. f(x) \wedge former(x)$	N $\lambda x. municipalities(x)$	$N \setminus N/NP$ $\lambda f \lambda x \lambda y. f(y) \wedge in(y, x)$	NP $Brandenburg$
\xrightarrow{N} $\lambda x. former(x) \wedge municipalities(x)$		$\xrightarrow{N \setminus N}$ $\lambda f \lambda y. f(y) \wedge in(y, Brandenburg)$	
\xleftarrow{N} $l_0 = \lambda x. former(x) \wedge municipalities(x) \wedge in(x, Brandenburg)$			

(b) Constant matches replace underspecified constants with Freebase concepts

$$l_0 = \lambda x. former(x) \wedge municipalities(x) \wedge in(x, Brandenburg)$$

$$l_1 = \lambda x. former(x) \wedge municipalities(x) \wedge in(x, Brandenburg)$$

$$l_2 = \lambda x. former(x) \wedge municipalities(x) \wedge location.containedby(x, Brandenburg)$$

$$l_3 = \lambda x. former(x) \wedge OpenRel(x, Municipality) \wedge location.containedby(x, Brandenburg)$$

$$l_4 = \lambda x. OpenType(x) \wedge OpenRel(x, Municipality) \wedge location.containedby(x, Brandenburg)$$

QA from Open IE

(a) CCG parse builds an underspecified semantic representation of the sentence.

Former	municipalities	in	Brandenburg
N/N $\lambda f \lambda x. f(x) \wedge former(x)$	N $\lambda x. municipalities(x)$	$N \setminus N/NP$ $\lambda f \lambda x \lambda y. f(y) \wedge in(y, x)$	NP $Brandenburg$
\xrightarrow{N} $\lambda x. former(x) \wedge municipalities(x)$		$\xrightarrow{N \setminus N}$ $\lambda f \lambda y. f(y) \wedge in(y, Brandenburg)$	
\xleftarrow{N} $l_0 = \lambda x. former(x) \wedge municipalities(x) \wedge in(x, Brandenburg)$			

(b) Constant matches replace underspecified constants with Freebase concepts

$$l_0 = \lambda x. former(x) \wedge municipalities(x) \wedge in(x, Brandenburg)$$

$$l_1 = \lambda x. former(x) \wedge municipalities(x) \wedge in(x, Brandenburg)$$

$$l_2 = \lambda x. former(x) \wedge municipalities(x) \wedge location.containedby(x, Brandenburg)$$

$$l_3 = \lambda x. former(x) \wedge OpenRel(x, Municipality) \wedge location.containedby(x, Brandenburg)$$

$$l_4 = \lambda x. OpenType(x) \wedge OpenRel(x, Municipality) \wedge location.containedby(x, Brandenburg)$$

- ▶ Why use the KB at all? Why not answer questions directly from text?
Like information retrieval!

Choi et al. (2015)

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
 - ▶ *What is the meaning of life?*
 - ▶ *What is 4+5?*
 - ▶ *What is the translation of [sentence] into French?* [McCann et al., 2018]

What can't KB QA systems do?

What can't KB QA systems do?

- ▶ What were the main causes of World War II? — requires summarization

What can't KB QA systems do?

- ▶ What were the main causes of World War II? — requires summarization
- ▶ Can you get the flu from a flu shot? — want IR to provide an explanation of the answer

What can't KB QA systems do?

- ▶ What were the main causes of World War II? — requires summarization
- ▶ Can you get the flu from a flu shot? — want IR to provide an explanation of the answer
- ▶ What temperature should I cook chicken to? — could be written down in a KB but probably isn't

What can't KB QA systems do?

- ▶ What were the main causes of World War II? — requires summarization
- ▶ Can you get the flu from a flu shot? — want IR to provide an explanation of the answer
- ▶ What temperature should I cook chicken to? — could be written down in a KB but probably isn't
- ▶ Today: can we do QA when it requires retrieving the answer from a passage?

Reading Comprehension

- ▶ “AI challenge problem”: answer question given context

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

Reading Comprehension

- ▶ “AI challenge problem”: answer question given context
- ▶ Recognizing Textual Entailment (2006)

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

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

Baselines

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.

- 2) What did James pull off of the shelves in the grocery store?
 - A) pudding
 - B) fries
 - C) food
 - D) splinters

Baselines

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.

2) What did James pull off of the shelves in the grocery store?

- A) pudding
- B) fries
- C) food
- D) splinters

Baselines

- ▶ N-gram matching: append question + each answer, return answer which gives highest n-gram overlap with a sentence

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.

- 2) What did James pull off of the shelves in the grocery store?
- A) pudding
 - B) fries
 - C) food
 - D) splinters

Baselines

- ▶ N-gram matching: append question + each answer, return answer which gives highest n-gram overlap with a sentence
- ▶ Parsing: find direct object of “pulled” in the document where the subject is James

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.

- 2) What did James pull off of the shelves in the grocery store?
- A) pudding
B) fries
C) food
D) splinters

Baselines

- ▶ N-gram matching: append question + each answer, return answer which gives highest n-gram overlap with a sentence
- ▶ Parsing: find direct object of “pulled” in the document where the subject is James
- ▶ Don’t need any complex semantic representations

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.

- 2) What did James pull off of the shelves in the grocery store?
- A) pudding
B) fries
C) food
D) splinters

Reading Comprehension

ngram sliding
window

	MC160 Test	MC500 Test
Baseline (SW+D)	66.25	56.67
RTE	59.79 [‡]	53.52
Combined	67.60	60.83 [‡]

- ▶ Classic textual entailment systems don't work as well as n-grams

Reading Comprehension

ngram sliding
window

	MC160 Test	MC500 Test
Baseline (SW+D)	66.25	56.67
RTE	59.79 [‡]	53.52
Combined	67.60	60.83 [‡]

- ▶ Classic textual entailment systems don't work as well as n-grams
- ▶ Scores are low partially due to questions spanning multiple sentences

Reading Comprehension

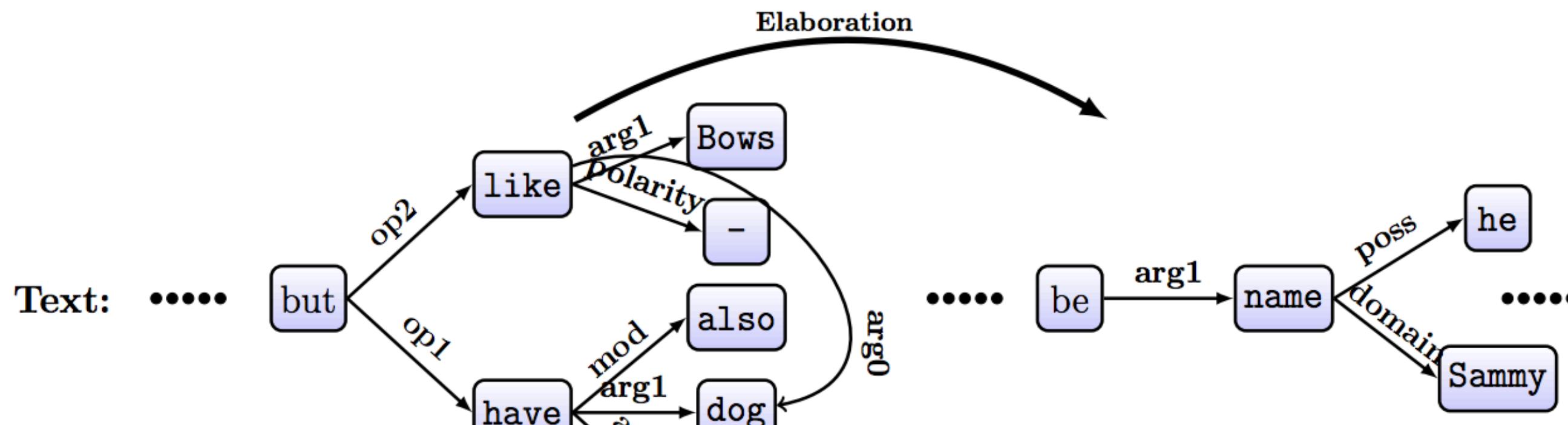
ngram sliding
window

	MC160 Test	MC500 Test
Baseline (SW+D)	66.25	56.67
RTE	59.79 [‡]	53.52
Combined	67.60	60.83 [‡]

- ▶ Classic textual entailment systems don't work as well as n-grams
- ▶ Scores are low partially due to questions spanning multiple sentences
- ▶ Unfortunately not much data to train better methods on (2000 questions)

MCTest: Better Systems

Text: ... Katie also has a dog, but he does not like Bows. ... His name is Sammy. ...

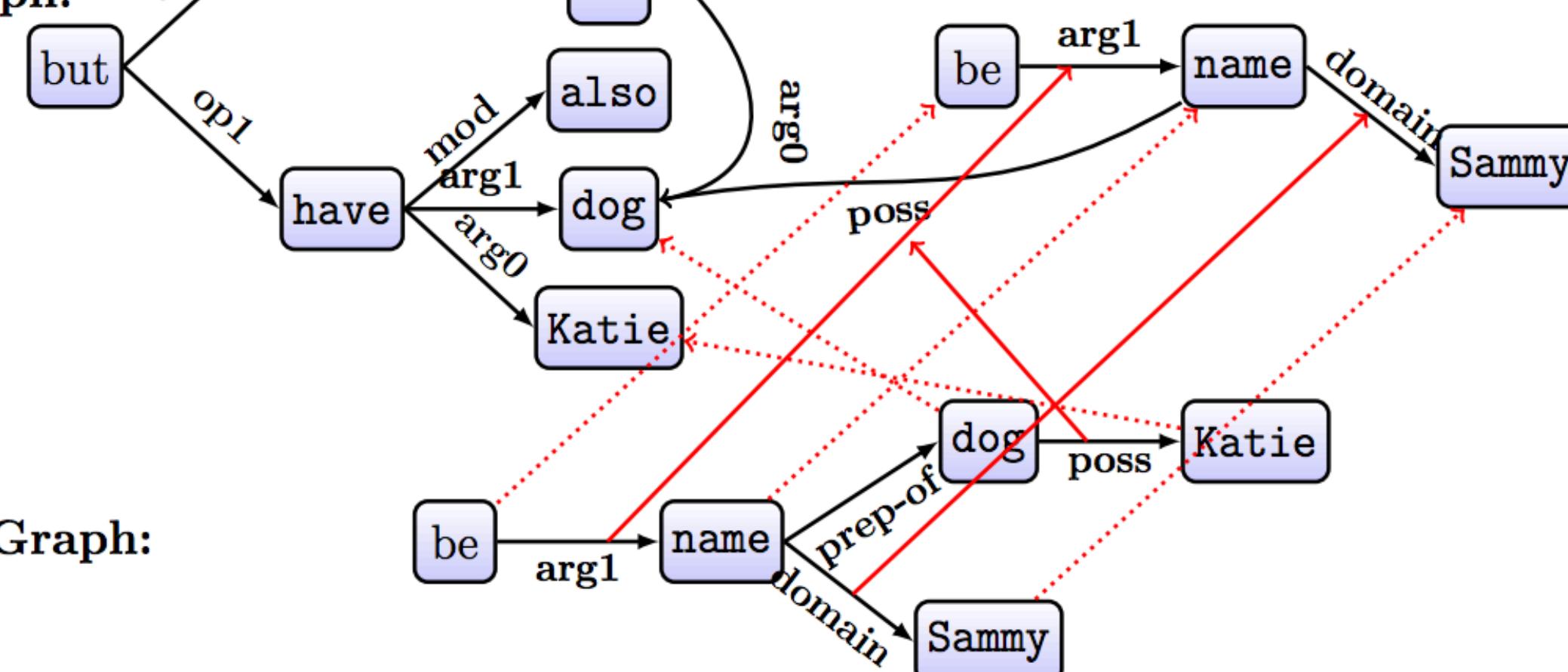


Text:

Snippet Graph:

Alignments:

Hypothesis Graph:



Hypothesis: Sammy is the name of Katie's dog.
Question: What is the name of Katie's dog. Answer: Sammy

- ▶ Match an AMR (abstract meaning representation) of the question against the original text
- ▶ 70% accuracy (roughly 10% better than baseline)

Sachan and Xing (2016)

Dataset Explosion

- ▶ 30+ QA datasets released since 2015
 - ▶ Children's Book Test, CNN/Daily Mail, SQuAD, TriviaQA, SearchQA, MS Marco, RACE, WikiHop, ...

Dataset Explosion

- ▶ 30+ QA datasets released since 2015
 - ▶ Children's Book Test, CNN/Daily Mail, SQuAD, TriviaQA, SearchQA, MS Marco, RACE, WikiHop, ...
- ▶ Question answering: questions are in natural language
 - ▶ Answers: multiple choice or require picking from the passage
 - ▶ Require human annotation

Dataset Explosion

- ▶ 30+ QA datasets released since 2015
 - ▶ Children’s Book Test, CNN/Daily Mail, SQuAD, TriviaQA, 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

Dataset Properties

- ▶ Axis 1: cloze task (fill in blank) vs. multiple choice vs. span-based vs. freeform generation
- ▶ Axis 2: what's the input?
 - ▶ One paragraph? One document? All of Wikipedia?
 - ▶ Some explicitly require linking between multiple sentences (MCTest, WikiHop, HotpotQA)
- ▶ Axis 3: what capabilities are needed to answer questions?
 - ▶ Finding simple information? Combining information across multiple sources? Commonsense knowledge?

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 it is hard to corner him."

"Are the boys big ?" queried Esther anxiously.

"Yes. Thirteen and fourteen and big for their age. You can't whip 'em -- that is the trouble. A man might, but they'd twist you around their fingers. You'll have your hands full, I'm afraid. But maybe they'll behave all right after all."

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 Mr. Baxter had exaggerated matters a little.

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 ? ''
8 queried Esther anxiously .
9 `` Yes .
10 Thirteen and fourteen and big for their age .
11 You ca n't whip 'em -- that is the trouble .
12 A man might , but they 'd twist you around their fingers .
13 You 'll have your hands full , I 'm afraid .
14 But maybe they 'll behave all right after all . ''
15 Mr. Baxter privately had no hope that they would , but Esther hoped for the best .
16 She could not believe that Mr. Cropper would carry his prejudices into a personal application .
17 This conviction was strengthened when he overtook her walking from school the next day and drove her home .
18 He was a big , handsome man with a very suave , polite manner .
19 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 .
20 Esther felt relieved .

Q: She thought that Mr. _____ had exaggerated matters a little .

C: Baxter, Cropper, Esther, course, fingers, manner, objection, opinion, right, spite.

a: Baxter

- ▶ 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)

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
rascals 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)

LAMBADA

Context: They tuned, discussed for a moment, then struck up a lively jig. Everyone joined in, turning the courtyard into an even more chaotic scene, people now dancing in circles, swinging and spinning in circles, everyone making up their own dance steps. I felt my feet tapping, my body wanting to move.

Target sentence: Aside from writing, I 've always loved _____.

Target word: dancing

- ▶ GPT/BERT can in general do very well at cloze tasks because this is what they're trained to do
- ▶ Hard to come up with plausible alternatives: “cooking”, “drawing”, “soccer”, etc. don't work in the above context

SWAG

- ▶ Dataset was constructed to be difficult for ELMo
- ▶ BERT subsequently got 20+% accuracy improvements and achieved human-level performance
- ▶ Problem: distractors too easy

The person blows the leaves from a grass area using the blower. The blower...

- a) puts the trimming product over her face in another section.
- b) is seen up close with different attachments and settings featured.
- c) continues to blow mulch all over the yard several times.
- d) blows beside them on the grass.

Span-Based Question Answering

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

One of the most famous people born in Warsaw was Maria Skłodowska-Curie, who achieved international recognition for her research on radioactivity and was the first female recipient of the Nobel Prize. Famous musicians include Władysław Szpilman and Frédéric Chopin. Though Chopin was born in the village of Żelazowa Wola, about 60 km (37 mi) from Warsaw, he moved to the city with his family when he was seven months old. Casimir Pulaski, a Polish general and hero of the American Revolutionary War, was born here in 1745.

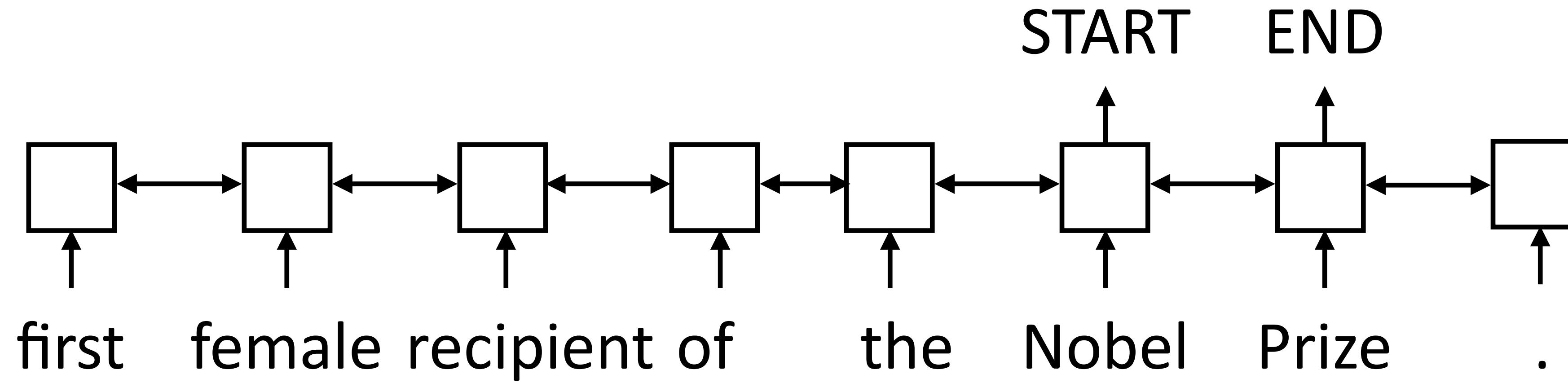
What was Maria Curie the first female recipient of?
Ground Truth Answers: Nobel Prize Nobel Prize Nobel Prize

What year was Casimir Pulaski born in Warsaw?
Ground Truth Answers: 1745 1745 1745

Who was one of the most famous people born in Warsaw?
Ground Truth Answers: Maria Skłodowska-Curie Maria Skłodowska-Curie Maria Skłodowska-Curie

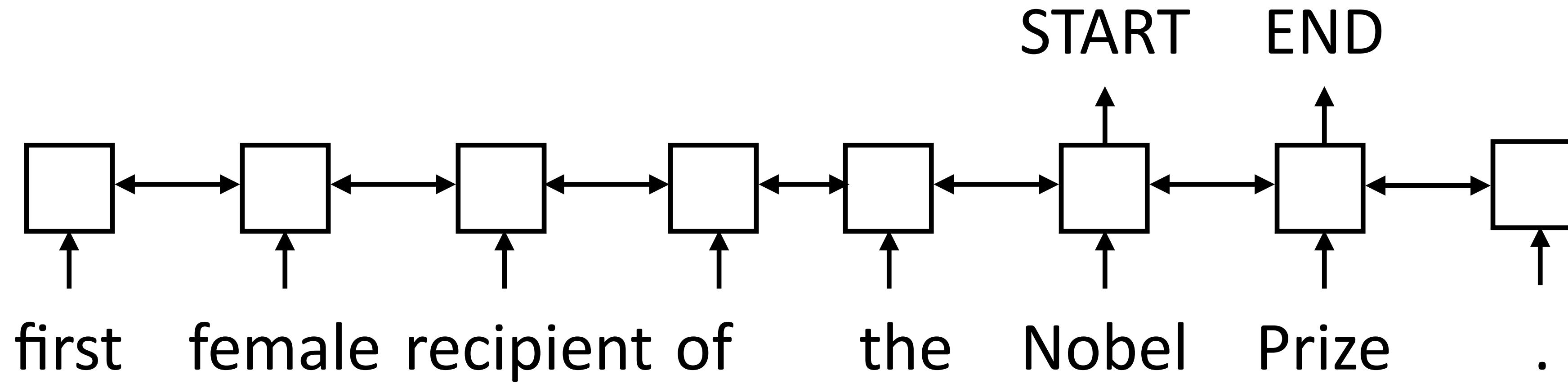
SQuAD

What was Marie Curie the first female recipient of?



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

Architectures

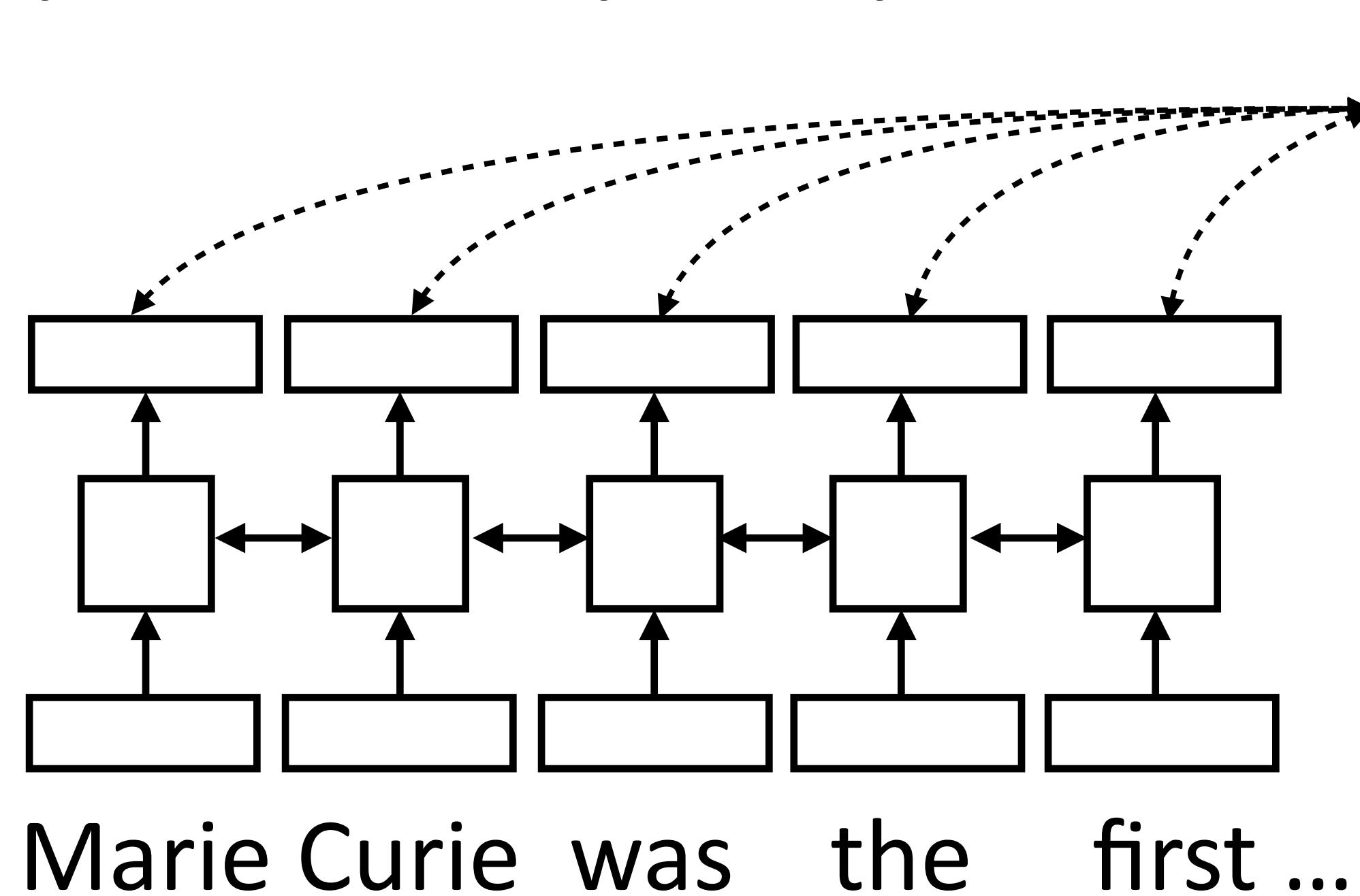
- ▶ Predict distributions over start and end points of the answer

$P(\text{end} \mid q, p)$ computed similarly

$$P(\text{start} = i \mid q, p) = \text{softmax}(p_i^\top W q)$$

encoding
of passage

BiLSTM
encoder



Who was the first
female recipient of
the Nobel Prize?

Training and Inference

- ▶ Train on labeled data with start and end points, maximize likelihood of correct decisions: $\log \sum_{i \in \text{gold starts}} p(\text{start} = i | p, q) + \log \sum_{i \in \text{gold ends}} p(\text{end} = i | p, q)$

In September 1958, Bank of America launched a new product called **BankAmericard** in Fresno. After a troubled gestation during which its creator resigned, **BankAmericard** went on to become the first **successful credit card**; that is, a financial instrument that was usable across a large number of merchants and also allowed **cardholders** to revolve a balance (earlier financial products could do one or the other but not both). In 1976, **BankAmericard** was **renamed** and spun off into a separate company known today as Visa Inc.

What was the name of the first successful credit card?

- ▶ Inference: maximize $P(\text{start}) + P(\text{end})$ with the constraint that $(\text{start}, \text{end})$ isn't too big a span

What do these models do?

Question: who caught a 16-yard pass on this drive ?

Answer: devin funchess

START

there would be no more scoring in the third quarter , but early in the fourth , the broncos drove to the panthers 41-yard line . on the next play , ealy knocked the ball out of manning 's hand as he was winding up for a pass , and then recovered it for carolina on the 50-yard line . a 16-yard reception by devin funchess and a 12-yard run by stewart then set up gano 's 39-yard field goal , cutting the panthers deficit to one score at 16â€“10 . the next three drives of the game would end in punts .

END

there would be no more scoring in the third quarter , but early in the fourth , the broncos drove to the panthers 41-yard line . on the next play , ealy knocked the ball out of manning 's hand as he was winding up for a pass , and then recovered it for carolina on the 50-yard line . a 16-yard reception by devin funchess and a 12-yard run by stewart then set up gano 's 39-yard field goal , cutting the panthers deficit to one score at 16â€“10 . the next three drives of the game would end in punts .

What do these models do?

Question: how many victorians are non - religious ?

Answer: 20 %

START

about 61.1 % of victorians describe themselves as christian . roman catholics form the single largest religious group in the state with 26.7 % of the victorian population , followed by anglicans and members of the uniting church . buddhism is the state 's largest non - christian religion , with 168,637 members as of the most recent census . victoria is also home of 152,775 muslims and 45,150 jews . hinduism is the fastest growing religion . around 20 % of victorians claim no religion . amongst those who declare a religious affiliation , church attendance is low .

END

about 61.1 % of victorians describe themselves as christian . roman catholics form the single largest religious group in the state with 26.7 % of the victorian population , followed by anglicans and members of the uniting church . buddhism is the state 's largest non - christian religion , with 168,637 members as of the most recent census . victoria is also home of 152,775 muslims and 45,150 jews . hinduism is the fastest growing religion . around 20 % of victorians claim no religion . amongst those who declare a religious affiliation , church attendance is low .

Why did this take off?

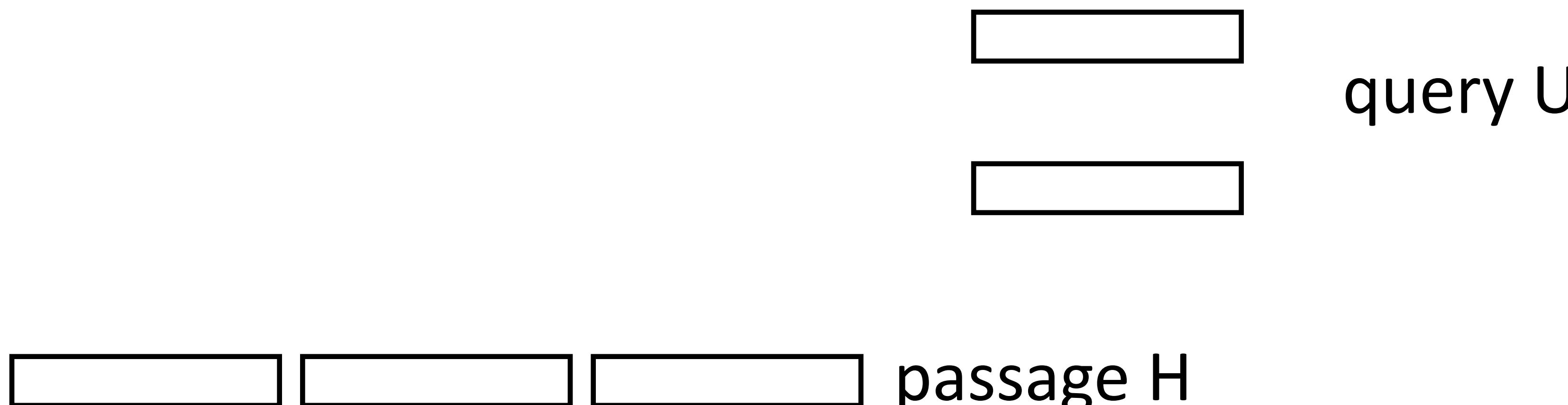
- ▶ SQuAD was **big**: >100,000 questions at a time when deep learning was exploding
- ▶ SQuAD was **pretty easy**: year-over-year progress for a few years until the dataset was essentially solved
- ▶ SQuAD had **room to improve**: ~50% performance from a logistic regression baseline (classifier with 180M features over constituents)

Bidirectional Attention Flow

- ▶ Passage (context) and query are both encoded with BiLSTMs

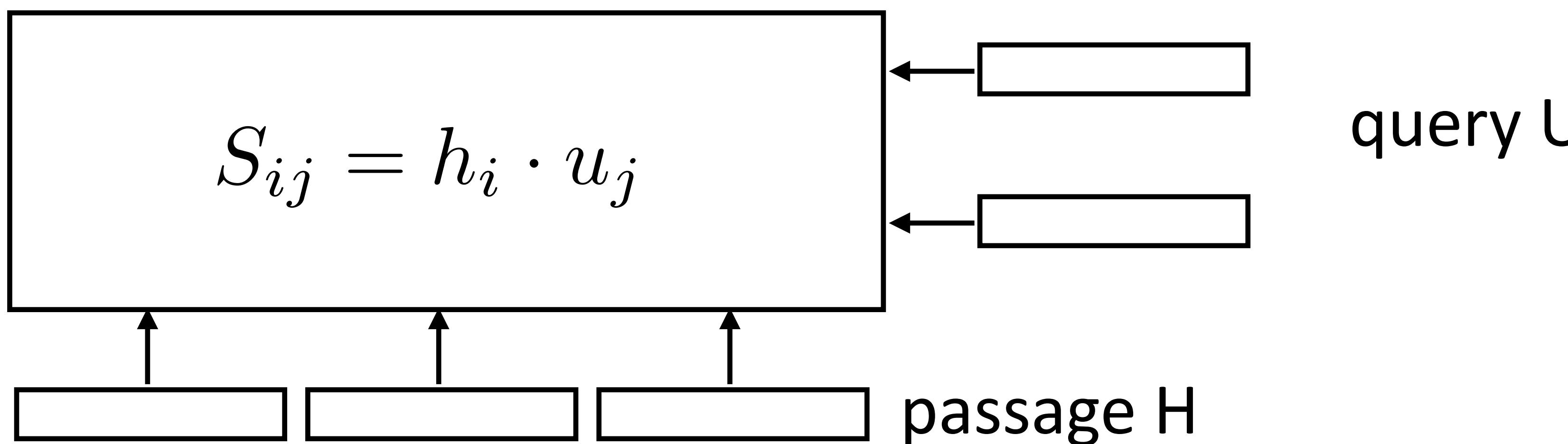
Bidirectional Attention Flow

- ▶ Passage (context) and query are both encoded with BiLSTMs



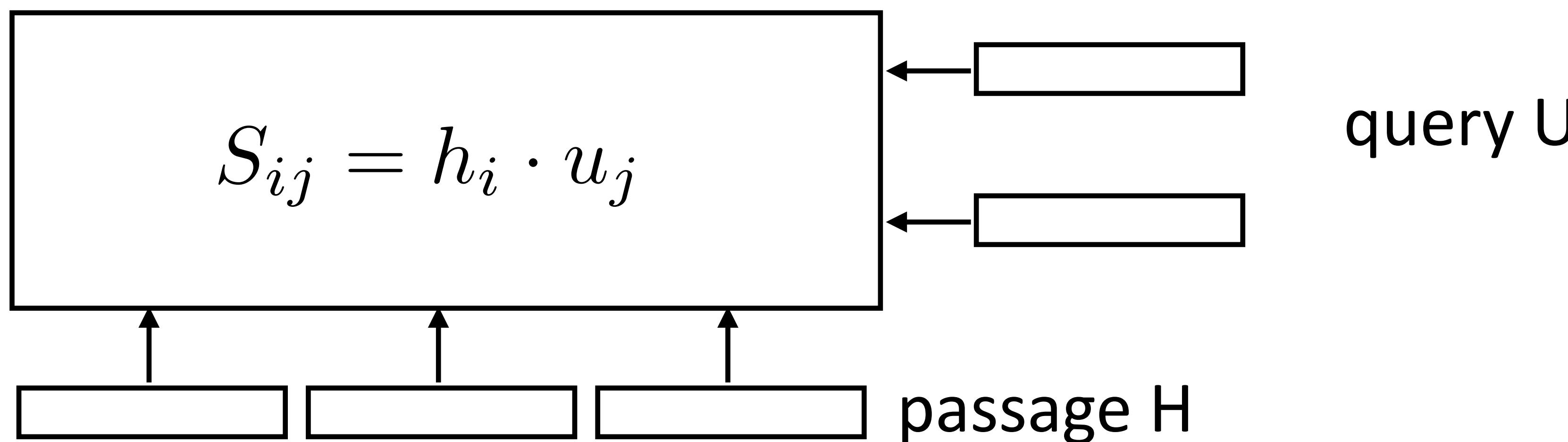
Bidirectional Attention Flow

- ▶ Passage (context) and query are both encoded with BiLSTMs



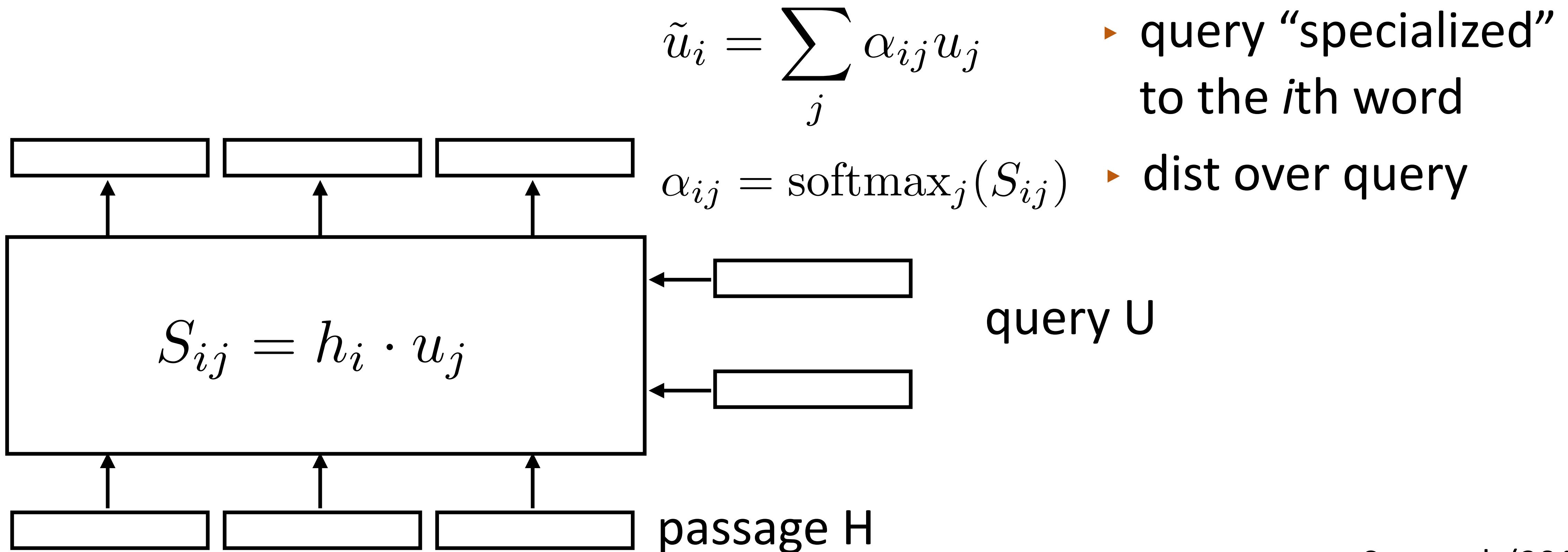
Bidirectional Attention Flow

- ▶ 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

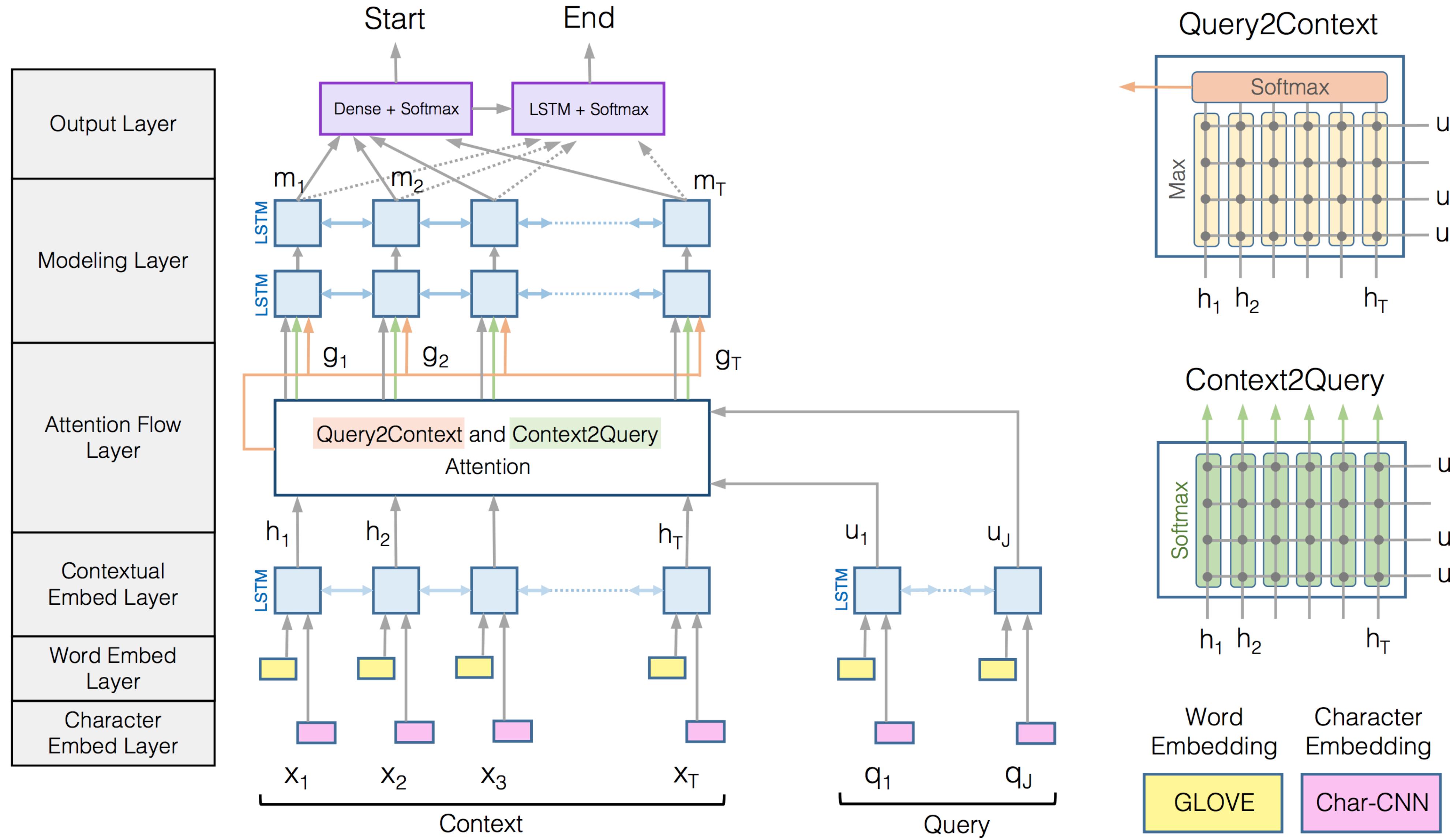


Bidirectional Attention Flow

- ▶ 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

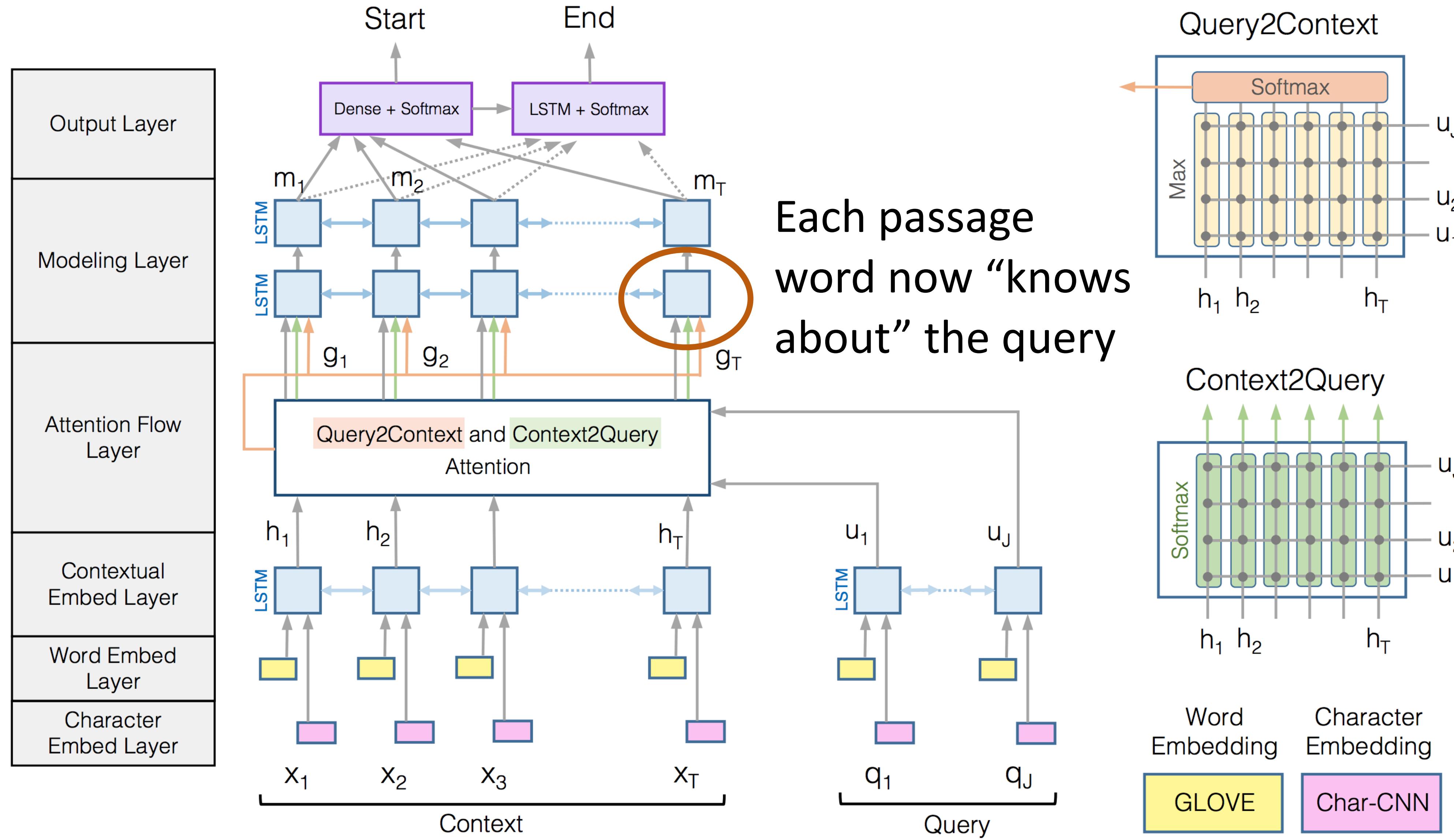


Bidirectional Attention Flow



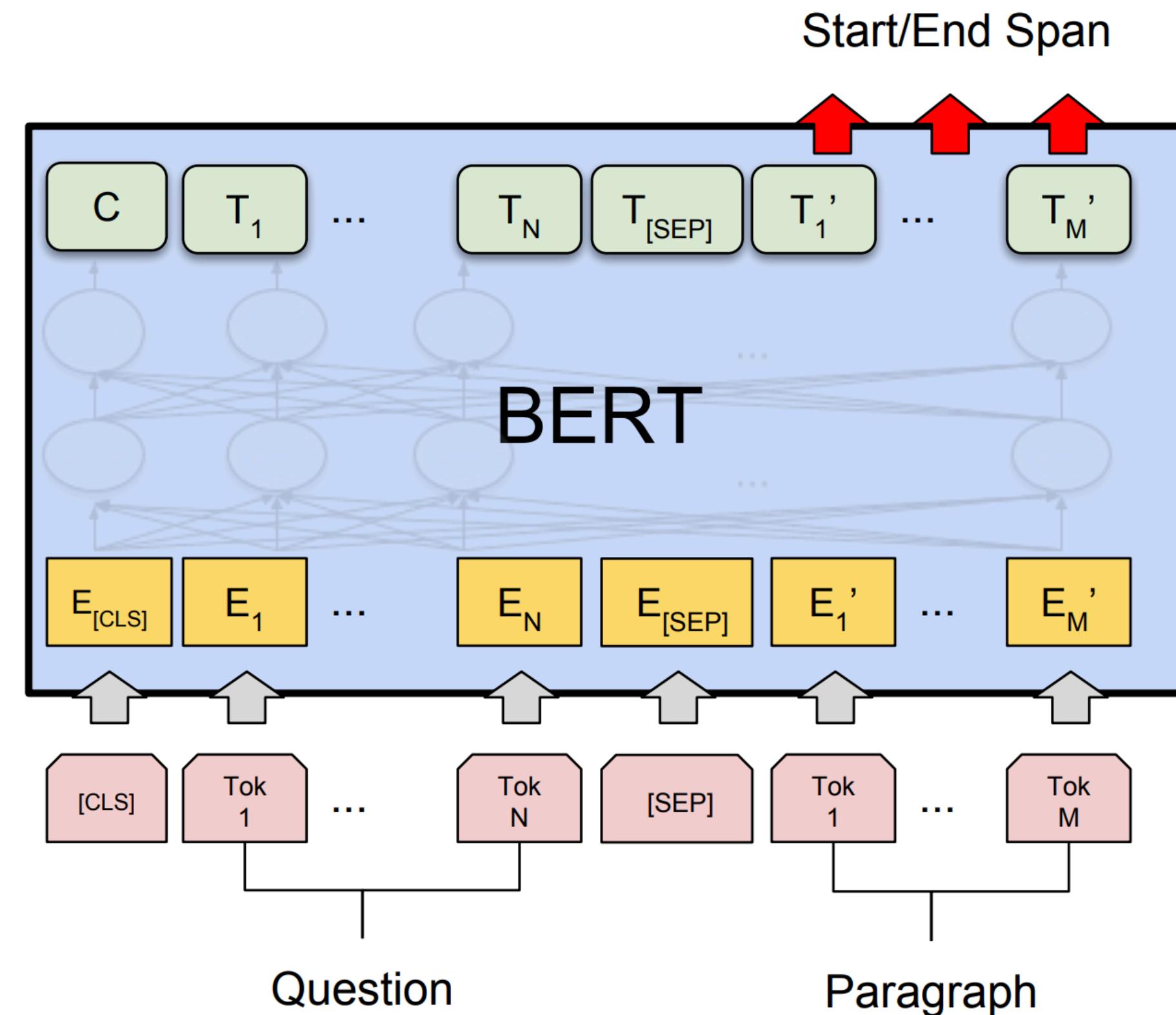
Seo et al. (2016)

Bidirectional Attention Flow



Seo et al. (2016)

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!

QA with BERT

► How does this work?

What was Marie Curie the first female recipient of ? [SEP] Marie Curie was the first female recipient of the Nobel Prize

SQuAD Results

SQuAD SOTA: Fall 18

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...)

SQuAD SOTA: Spring 19

Rank	Model	EM	F1	
	Human Performance <i>Stanford University</i> <i>(Rajpurkar & Jia et al. '18)</i>	86.831	89.452	
1	BERT + DAE + AoA (ensemble) <i>Joint Laboratory of HIT and iFLYTEK Research</i>	87.147	89.474	
2	Mar 20, 2019	BERT + ConvLSTM + MTL + Verifier (ensemble) <i>Layer 6 AI</i>	86.730	89.286
3	Mar 15, 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 SOTA: Fall 19

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

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

What are these models learning?

(Answer = Stanford University)

Question: Where did the Broncos practice for the Super Bowl ?

Passage: The Panthers used the San Jose State practice facility and stayed at the San Jose Marriott . The Broncos practiced at Stanford University and stayed at the Santa Clara Marriott .

(d) Erasure exact search optima.

Question: Where did the Broncos practice for the Super Bowl ?

Passage: The Panthers used the San Jose State practice facility and stayed at the San Jose Marriott . The Broncos practiced at Stanford University and stayed at the Santa Clara Marriott .

(a) Integrated Gradient ([Sundararajan et al., 2017](#)).

► Are these good explanations?

What are these models learning?

< s >
Who
did
the
Panthers
face
in
the
NFC
Championship
Game
?
</s>
The
Panthers
then
blew
out
the
Arizona
Cardinals
in
the
NFC
Championship
Game
,
49-15
,
racking
up
487
yards
and
forcing
seven
turnovers
. </s>

< s >
Who
did
the
Panthers
face
in
the
NFC
Championship
Game
?
</s>

Pairwise explanation method: explains predictions
in terms of associations between words

Ye, Nair, Durrett (2021)

What are these models learning?

<|>
What
typeface
are
the
letters
in
the
iconic
ABC
logo
reminiscent
of
?
</|>
Paul
Rand
redesigned
the
ABC
logo
into
its
best
-
known
(and
current
)
form
The
letters
are
strongly
reminiscent
of
the
Bauhaus
typeface

<|>
What
typeface
are
the
letters
in
the
iconic
ABC
logo
reminiscent
of
?
</|>

ABC isn't used at all! The model is mostly using
the fact that only one typeface is in the context

Ye, Nair, Durrett (2021)

Takeaways

- ▶ Lots of problems with current QA settings, lots of new datasets
- ▶ Models can often work well for one QA task but don't generalize
- ▶ We still don't have (solvable) QA settings which seem to require really complex reasoning as opposed to surface-level pattern recognition