

# Natural Language Processing with Deep Learning

## CS224N/Ling284

### Lecture 11: Question Answering

Danqi Chen  
Princeton University

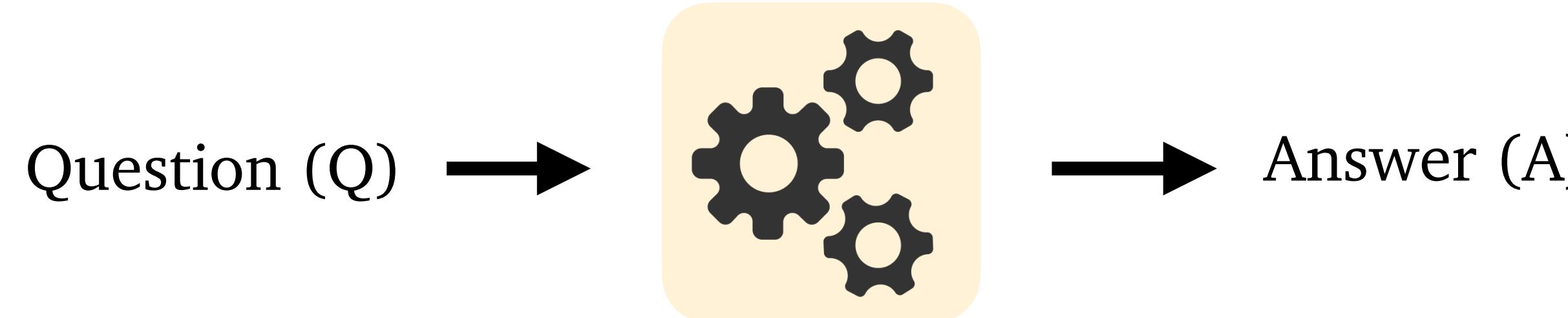
# Lecture plan

1. What is question answering? (10 mins)
2. Reading comprehension (50 mins)
  - ✓ How to answer questions over **a single passage of text**
3. Open-domain (textual) question answering (20 mins)
  - ✓ How to answer questions over **a large collection of documents**

Your default final project!

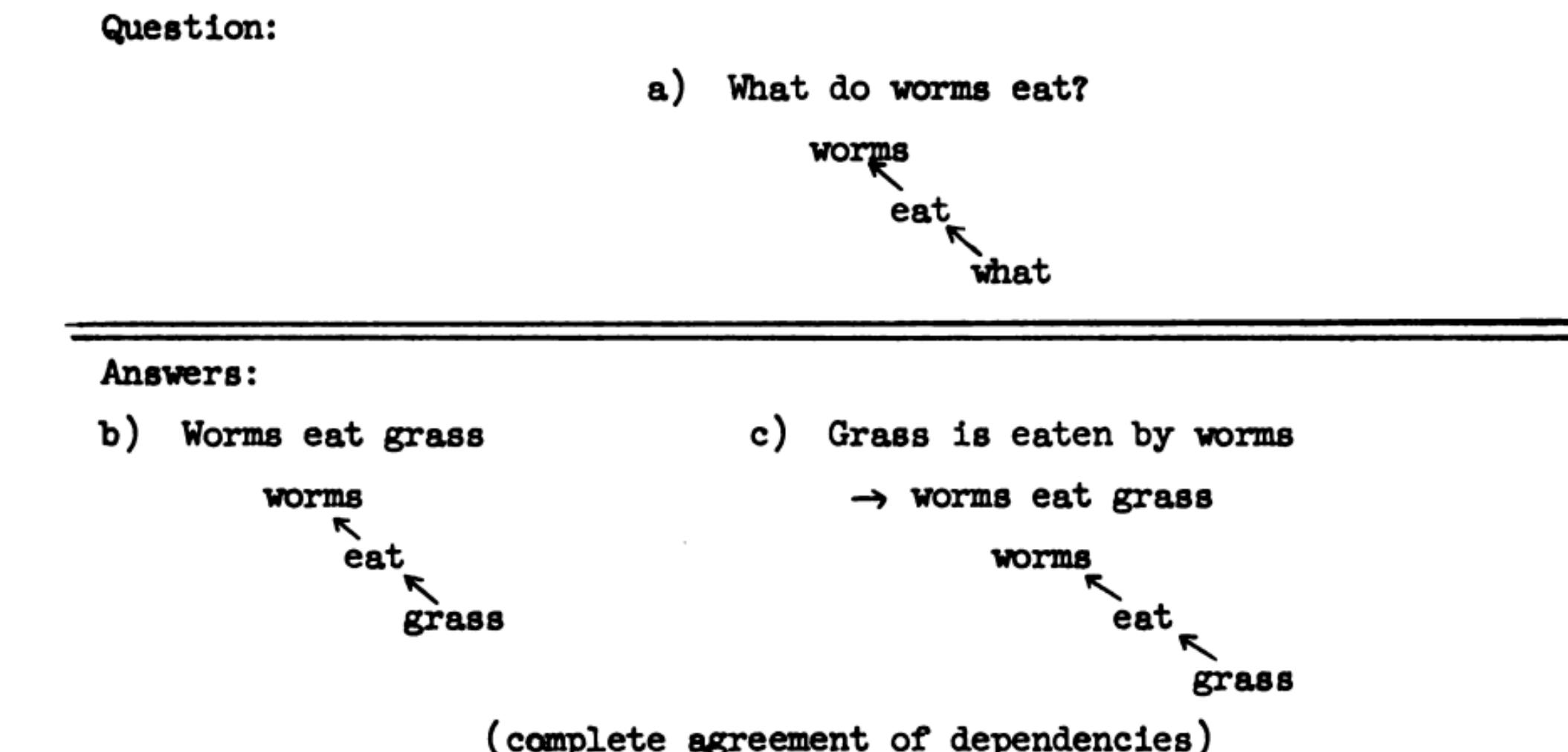


# 1. What is question answering?

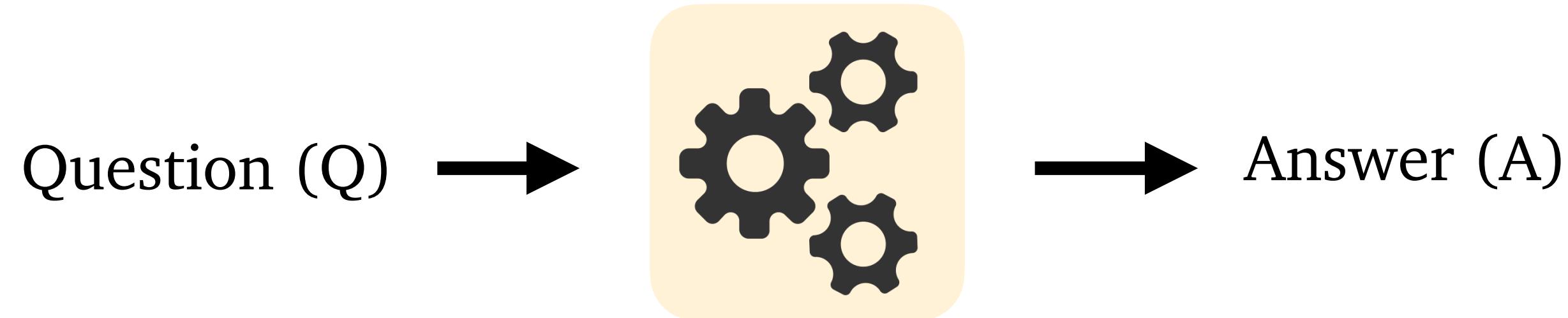


The goal of question answering is to build systems that **automatically** answer questions posed by humans in a **natural language**

The earliest QA systems  
dated back to 1960s!  
(Simmons et al., 1964)



# Question answering: a taxonomy



- What information source does a system build on?
  - 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/no, ...

# Lots of practical applications

A screenshot of a Google search results page. The search bar at the top contains the query "Where is the deepest lake in the world?". Below the search bar, there are several navigation links: "All" (which is highlighted in blue), "Maps", "Images", "News", "Videos", and "More". To the right of these are "Settings" and "Tools". A message indicates "About 21,100,000 results (0.71 seconds)". Below this, there are four thumbnail images: a landscape view of a lake, a 3D cross-section diagram of Lake Baikal, another landscape view of a lake, and a satellite map showing a deep blue lake.

## Siberia

Lake **Baikal**, in Siberia, holds the distinction of being both the deepest lake in the world and the largest freshwater lake, holding more than 20% of the unfrozen fresh water on the surface of Earth.

# Lots of practical applications

Google How can I protect myself from COVID-19? X |

All Images News Shopping Videos More Settings Tools

The best way to prevent illness is to avoid being exposed to this virus. Learn how COVID-19 spreads and practice these actions to help prevent the spread of this illness.

To help prevent the spread of COVID-19:

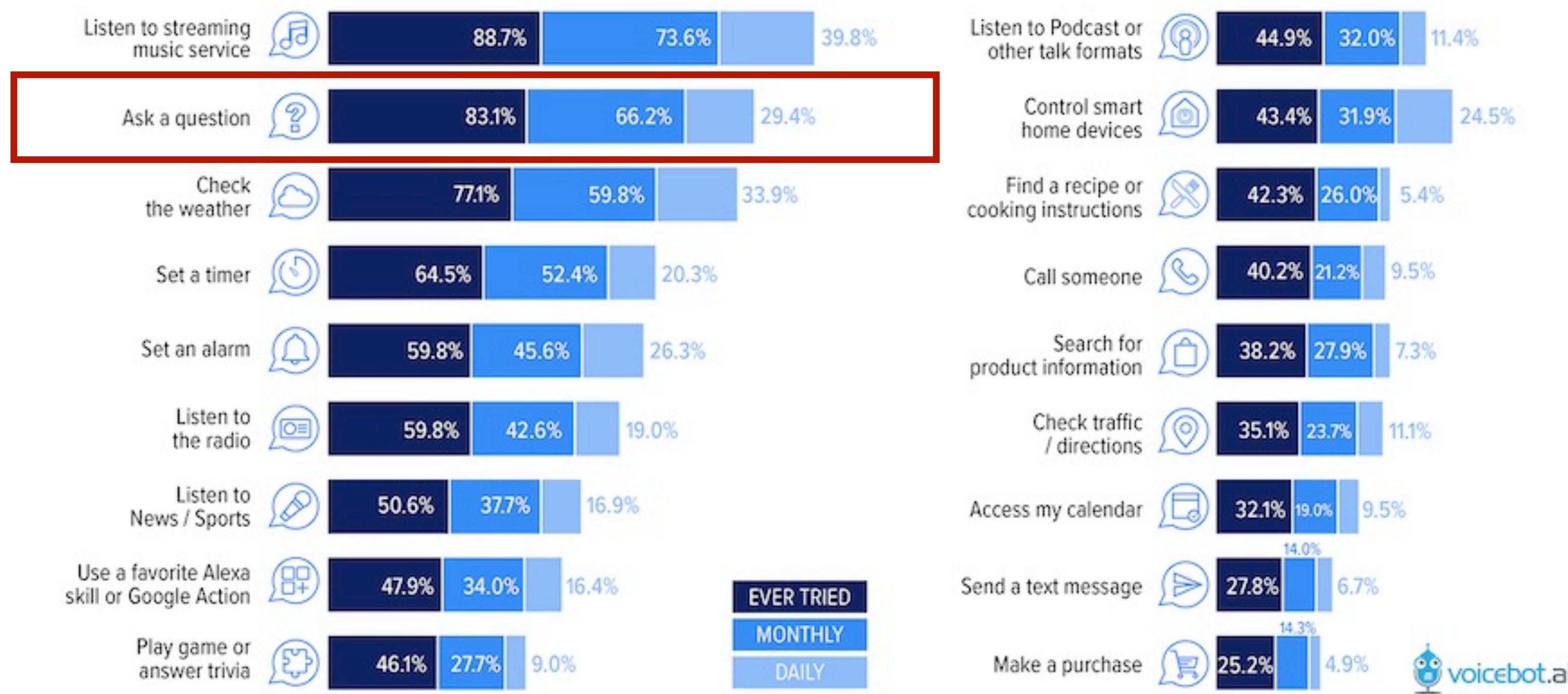
- Cover your mouth and nose with a mask when around people who don't live with you. Masks work best when everyone wears one.
- Stay at least 6 feet (about 2 arm lengths) from others.
- Avoid crowds. The more people you are in contact with, the more likely you are to be exposed to COVID-19.
- Avoid unventilated indoor spaces. If indoors, bring in fresh air by opening windows and doors.
- Clean your hands often, either with soap and water for 20 seconds or a hand sanitizer that contains at least 60% alcohol.
- Get vaccinated against COVID-19 when it's your turn.
- Avoid close contact with people who are sick.
- Cover your cough or sneeze with a tissue, then throw the tissue in the trash.
- Clean and disinfect frequently touched objects and surfaces daily.

Learn more on [cdc.gov](#)

For informational purposes only. Consult your local medical authority for advice.

# Lots of practical applications

Smart Speaker Use Case Frequency January 2020



Source: Voicebot.ai 2020

# IBM Watson beat Jeopardy champions



IBM Watson defeated two of Jeopardy's greatest champions in 2011

# IBM Watson beat Jeopardy champions

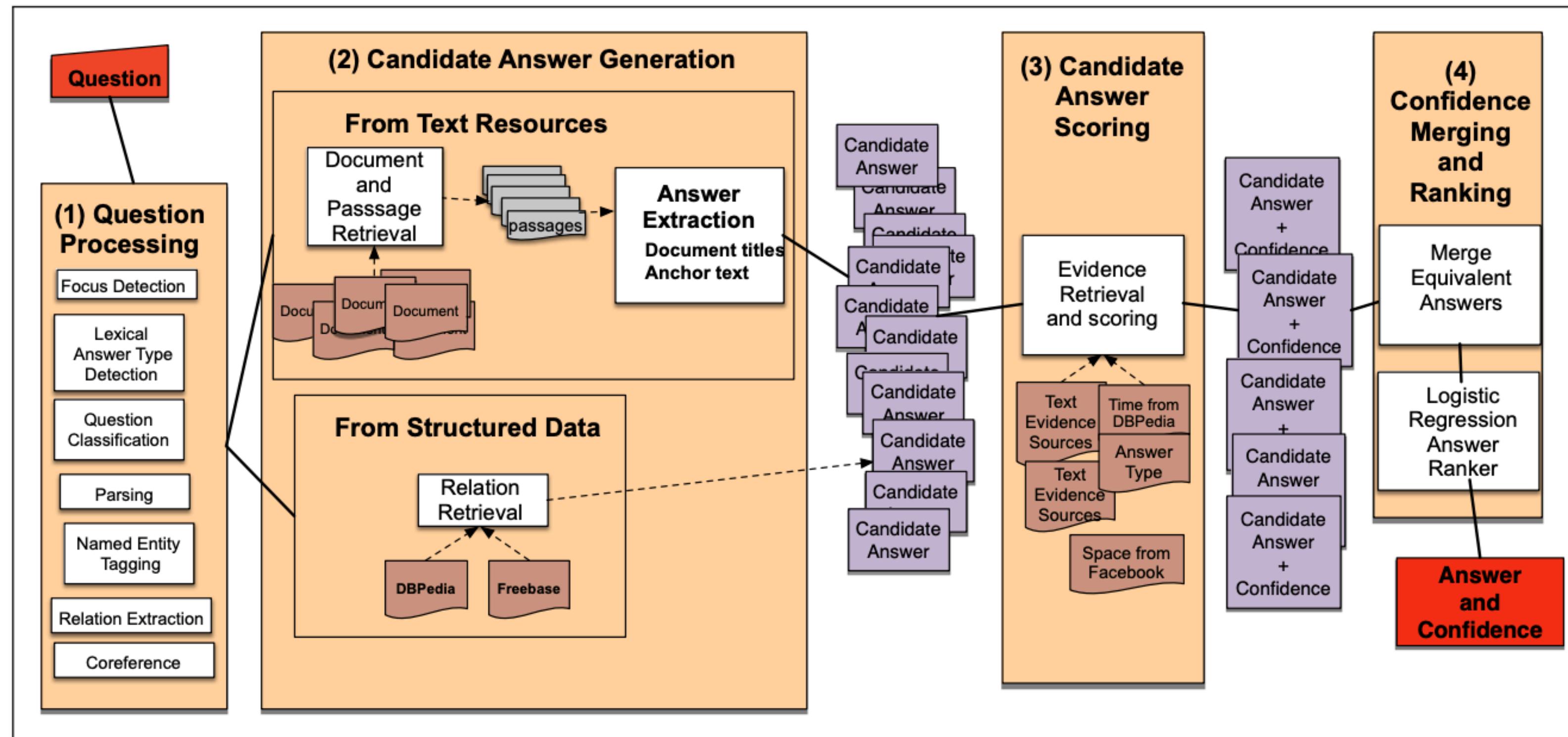


Image credit: J & M, edition 3

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

# Question answering in deep learning era

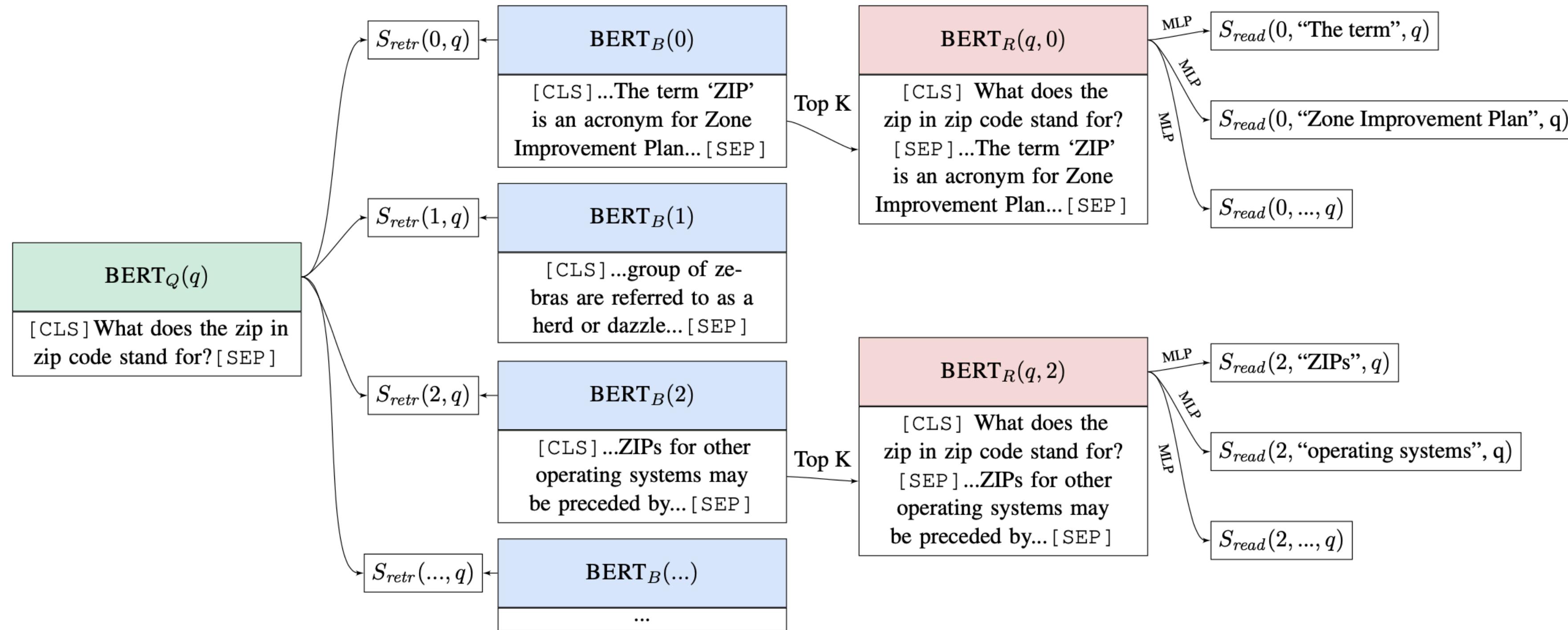


Image credit: (Lee et al., 2019)

Almost all the state-of-the-art question answering systems are built on top of end-to-end training and pre-trained language models (e.g., BERT)!

# Beyond textual QA problems

Today, we will mostly focus on how to answer questions based on **unstructured text**.

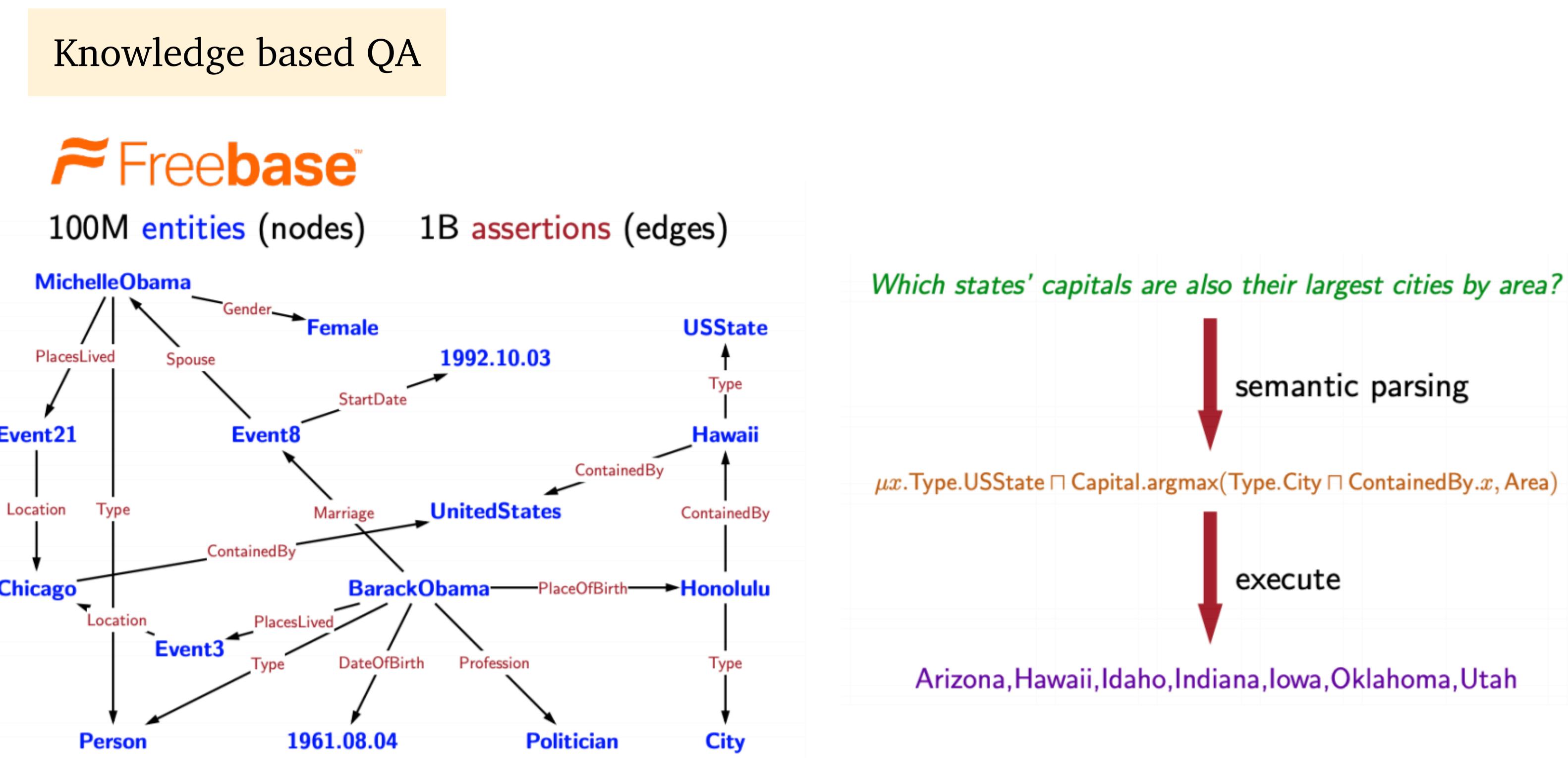


Image credit: Percy Liang

# Beyond textual QA problems

Today, we will mostly focus on how to answer questions based on **unstructured text**.

Visual QA



What color are her eyes?  
What is the mustache made of?



How many slices of pizza are there?  
Is this a vegetarian pizza?

(Antol et al., 2015): Visual Question Answering

## 2. Reading comprehension

**Reading comprehension** = comprehend a passage of text and answer questions about its content (P, Q) → A

Tesla was the fourth of five children. He had an older brother named Dane and three sisters, Milka, Angelina and Marica. Dane was killed in a horse-riding accident when Nikola was five. In 1861, Tesla attended the "Lower" or "Primary" School in Smiljan where he studied German, arithmetic, and religion. In 1862, the Tesla family moved to Gospić, Austrian Empire, where Tesla's father worked as a pastor. Nikola completed "Lower" or "Primary" School, followed by the "Lower Real Gymnasium" or "Normal School."

**Q:** What language did Tesla study while in school?

**A:** German

## 2. Reading comprehension

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

Kannada language is the official language of Karnataka and spoken as a native language by about 66.54% of the people as of 2011. Other linguistic minorities in the state were Urdu (10.83%), Telugu language (5.84%), Tamil language (3.45%), Marathi language (3.38%), Hindi (3.3%), Tulu language (2.61%), Konkani language (1.29%), Malayalam (1.27%) and Kodava Takk (0.18%). In 2007 the state had a birth rate of 2.2%, a death rate of 0.7%, an infant mortality rate of 5.5% and a maternal mortality rate of 0.2%. The total fertility rate was 2.2.

Q: Which linguistic minority is larger, Hindi or Malayalam?

A: Hindi

# Why do we care about this problem?

- Useful for many practical applications
- Reading comprehension is an important testbed for evaluating how well computer systems understand human language
  - Wendy Lehnert 1977: “Since questions can be devised to query **any aspect** of text comprehension, the ability to answer questions is the **strongest possible demonstration of understanding.**”
- Many other NLP tasks can be reduced to a reading comprehension problem:

## Information extraction

(Barack Obama, educated\_at, ?)

Question: Where did Barack Obama graduate from?

Passage: Obama was born in Honolulu, Hawaii.  
After graduating from Columbia University in 1983,  
he worked as a community organizer in Chicago.

(Levy et al., 2017)

## Semantic role labeling

UCD **finished** the 2006 championship as Dublin champions ,  
by **beating** St Vincents in the final .

Who finished something? - UCD

What did someone finish? - the 2006 championship

What did someone finish something as? - Dublin champions

How did someone finish something? - by beating St Vincents in the final

Who beat someone? - UCD

When did someone beat someone? - in the final

Who did someone beat? - St Vincents

(He et al., 2015)

# Stanford question answering dataset (SQuAD)

- 100k annotated (passage, question, answer) triples

Large-scale supervised datasets are also a key ingredient for training effective neural models for reading comprehension!

- 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!

- SQuAD still remains the most popular reading comprehension dataset; it is “almost solved” today and the state-of-the-art exceeds the estimated human performance.

---

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under **gravity**. The main forms of precipitation include drizzle, rain, sleet, snow, **graupel** and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals **within a cloud**. Short, intense periods of rain in scattered locations are called “showers”.

What causes precipitation to fall?  
**gravity**

What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?  
**graupel**

Where do water droplets collide with ice crystals to form precipitation?  
**within a cloud**

---

# Stanford question answering dataset (SQuAD)

- **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 (a, an, the, punctuations are removed) and take max scores. Finally, we take the average of all the examples for both exact match and F1.
- Estimated human performance: EM = 82.3, F1 = 91.2

Q: What did Tesla do in December 1878?

A: {left Graz, left Graz, left Graz and severed all relations with his family}

Prediction: {left Graz and served}

Exact match:  $\max\{0, 0, 0\} = 0$

F1:  $\max\{0.67, 0.67, 0.61\} = 0.67$

# Neural models for reading comprehension

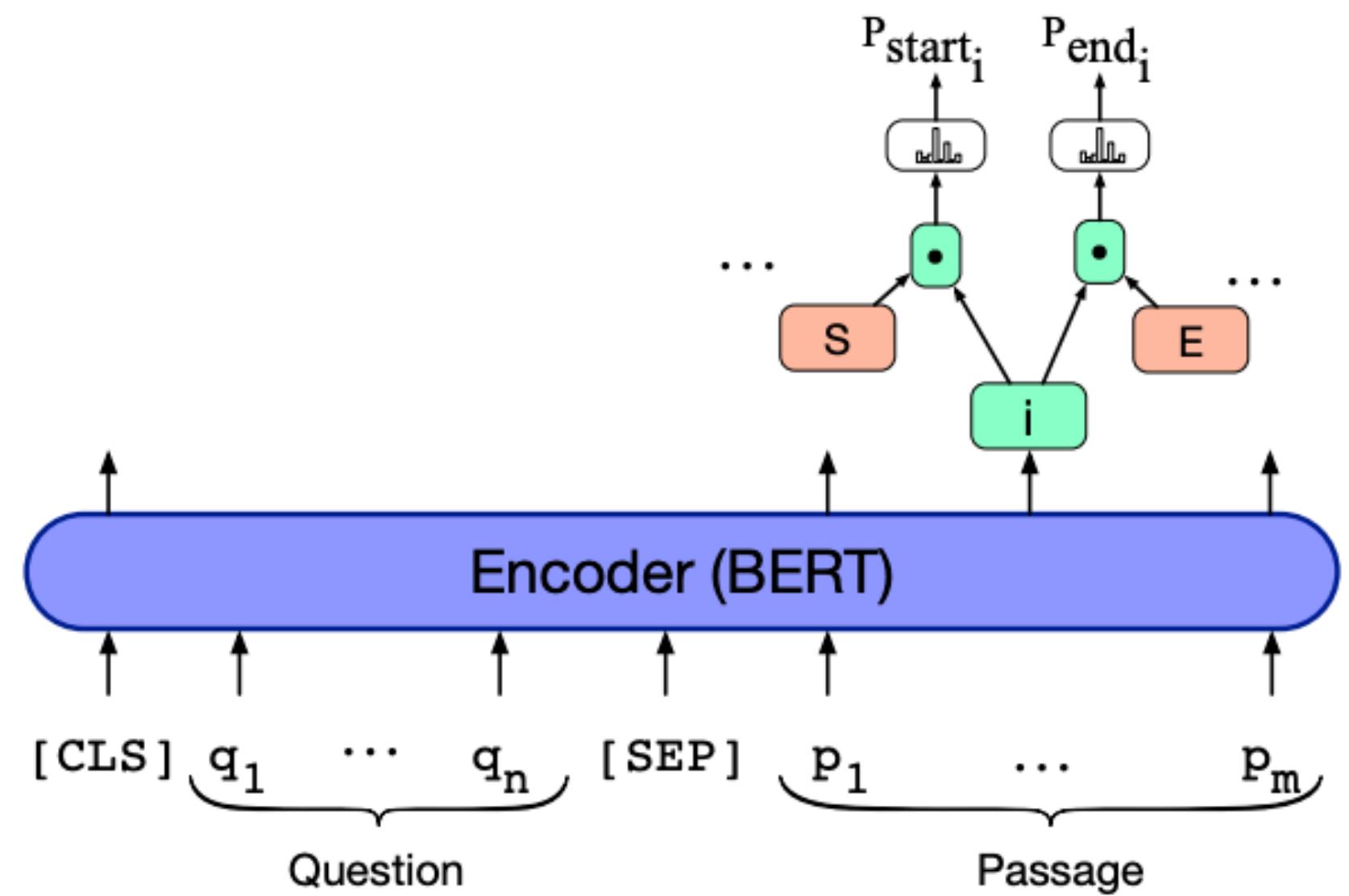
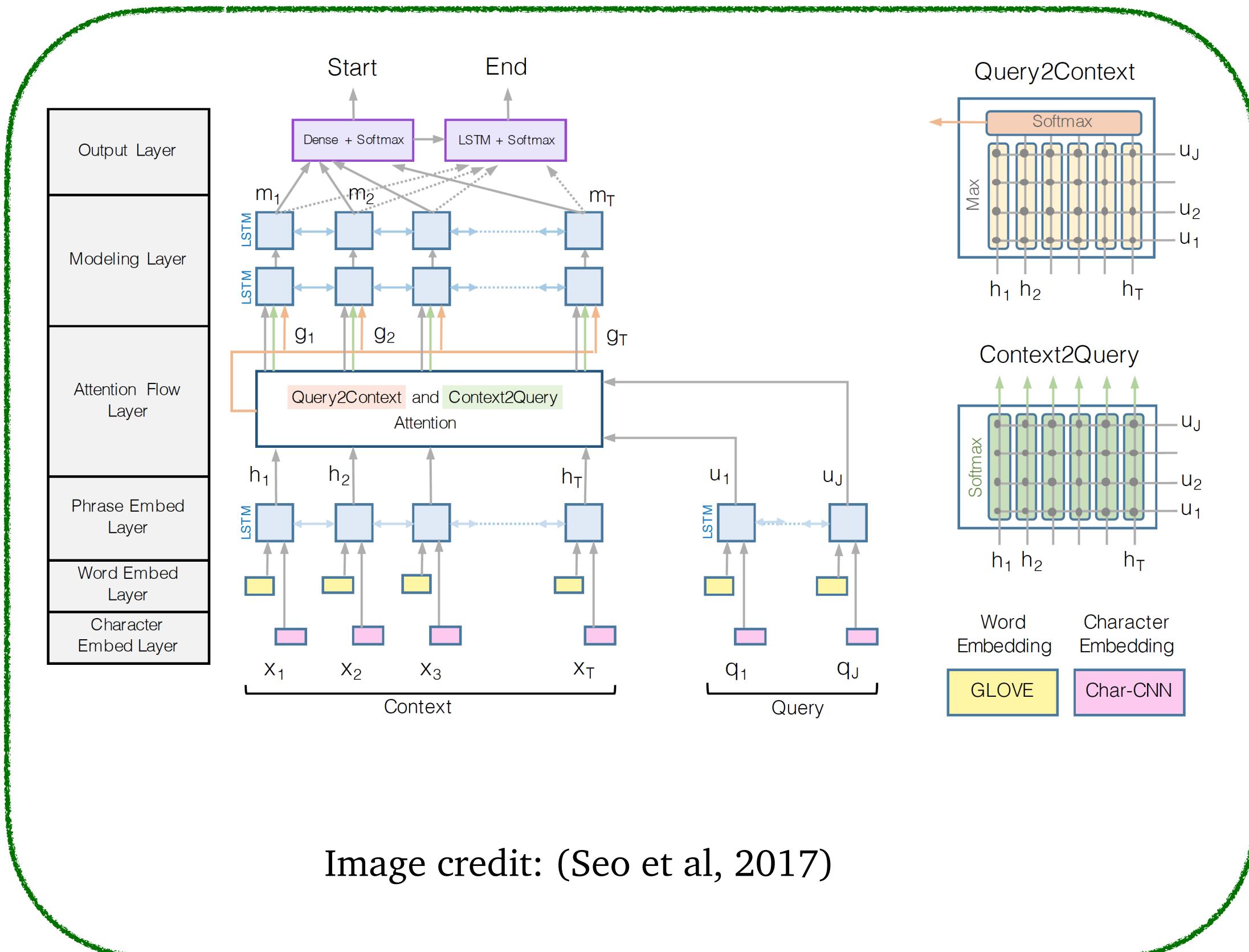
## How can we build a model to solve SQuAD?

(We are going to use **passage**, **paragraph** and **context**, as well as **question** and **query** interchangeably)

- Problem formulation
  - Input:  $C = (c_1, c_2, \dots, c_N), Q = (q_1, q_2, \dots, q_M), c_i, q_i \in V$  N \sim 100, M \sim 15
  - Output:  $1 \leq \text{start} \leq \text{end} \leq N$  answer is a span in the passage
- A family of LSTM-based models with attention (2016-2018)

Attentive Reader (Hermann et al., 2015), Stanford Attentive Reader (Chen et al., 2016), Match-LSTM (Wang et al., 2017), BiDFA (Seo et al., 2017), Dynamic coattention network (Xiong et al., 2017), DrQA (Chen et al., 2017), R-Net (Wang et al., 2017), ReasoNet (Shen et al., 2017)..
- Fine-tuning BERT-like models for reading comprehension (2019+)

# LSTM-based vs BERT models

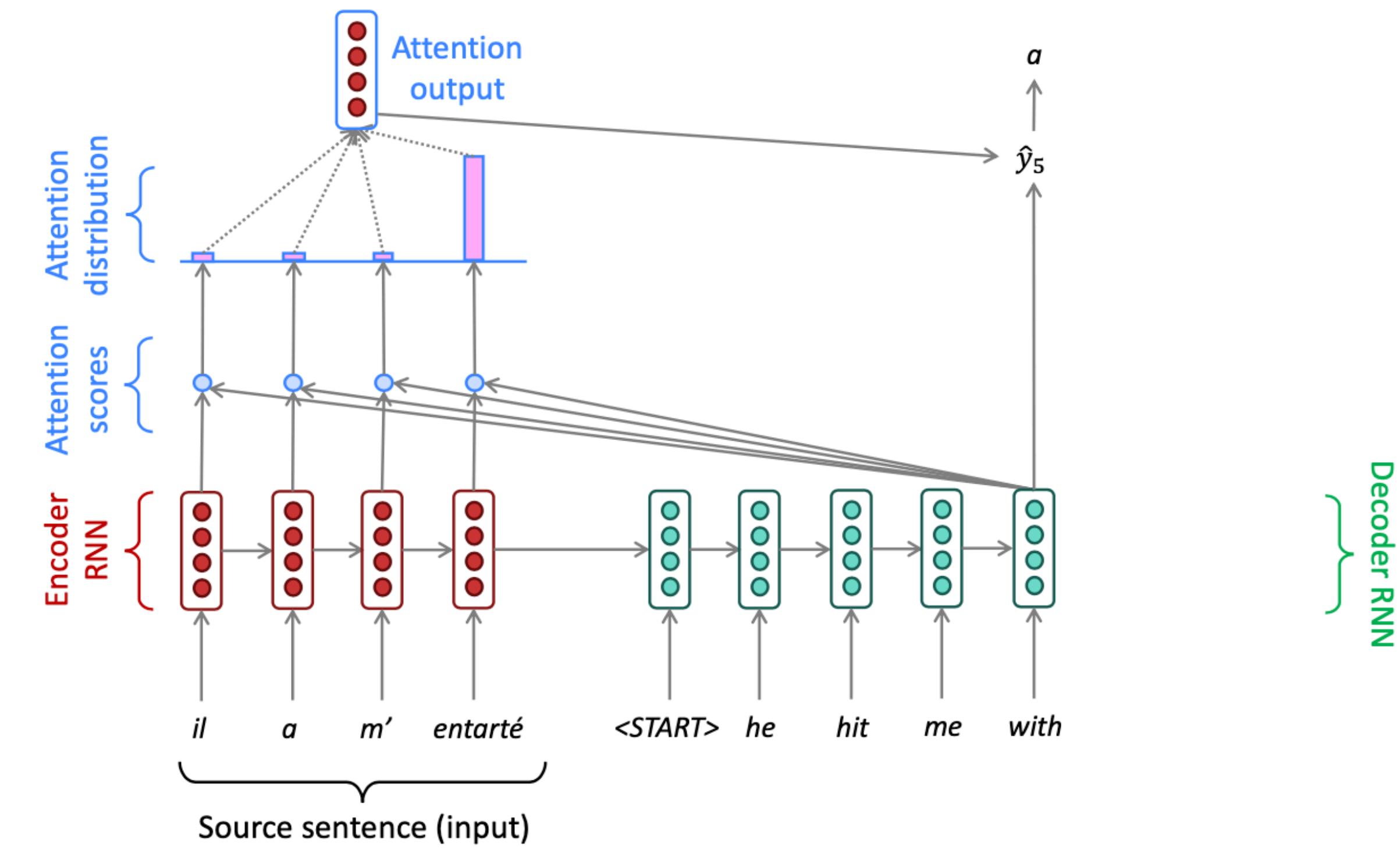


# Recap: seq2seq model with attention

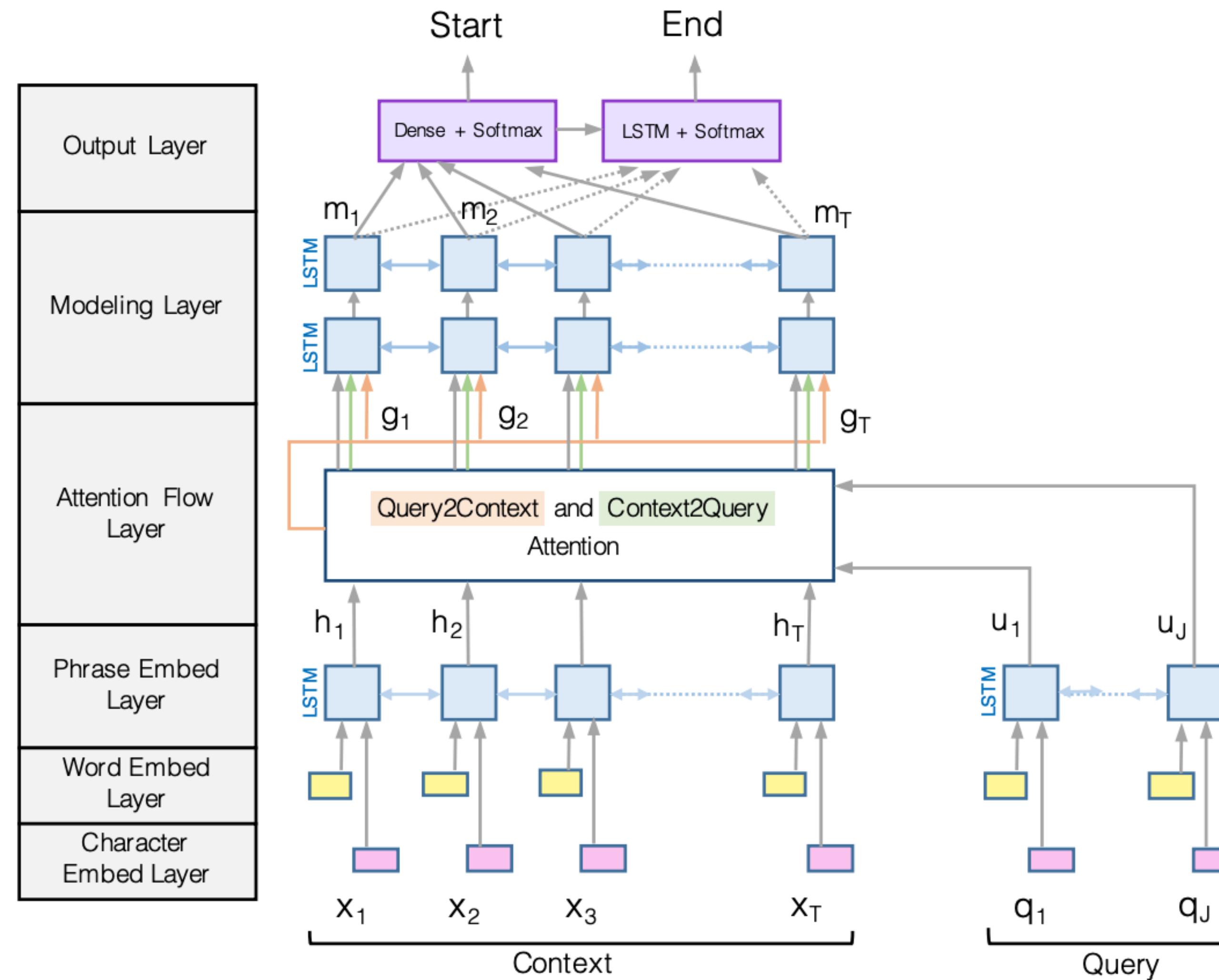
- Instead of source and target sentences, we also have two sequences: passage and question (lengths are imbalanced)
- We need to model which words in the passage are most relevant to the question (and which question words)

Attention is the key ingredient here, similar to which words in the source sentence are most relevant to the current target word...

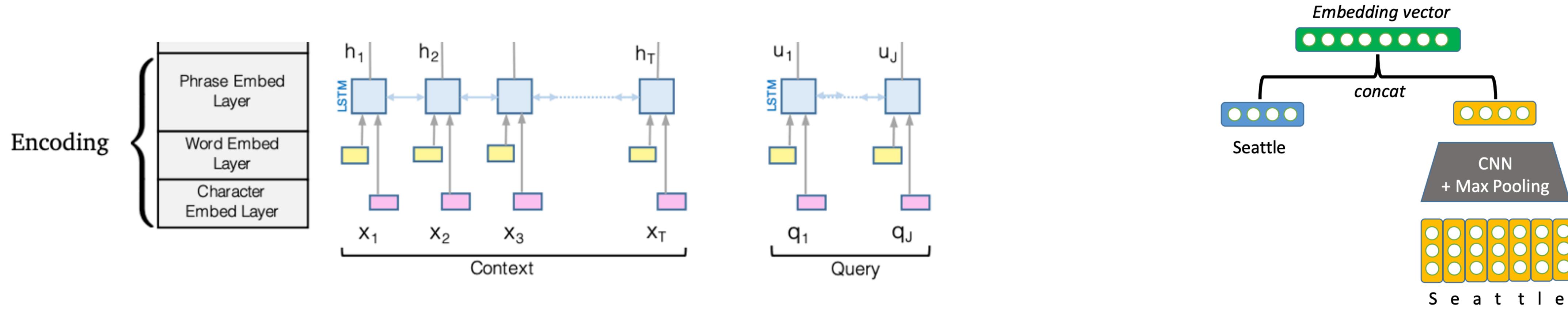
- We don't need an autoregressive decoder to generate the target sentence word-by-word. Instead, we just need to train two classifiers to predict the start and end positions of the answer!



# BiDAF: the Bidirectional Attention Flow model



# BiDAF: Encoding



- Use a concatenation of word embedding (GloVe) and character embedding (CNNs over character embeddings) for each word in context and query.

$$e(c_i) = f([\text{GloVe}(c_i); \text{charEmb}(c_i)])$$

$$e(q_i) = f([\text{GloVe}(q_i); \text{charEmb}(q_i)])$$

*f: highway networks omitted here*

- Then, use two bidirectional LSTMs separately to produce contextual embeddings for both context and query.

$$\vec{\mathbf{c}}_i = \text{LSTM}(\vec{\mathbf{c}}_{i-1}, e(c_i)) \in \mathbb{R}^H$$

$$\overleftarrow{\mathbf{c}}_i = \text{LSTM}(\overleftarrow{\mathbf{c}}_{i+1}, e(c_i)) \in \mathbb{R}^H$$

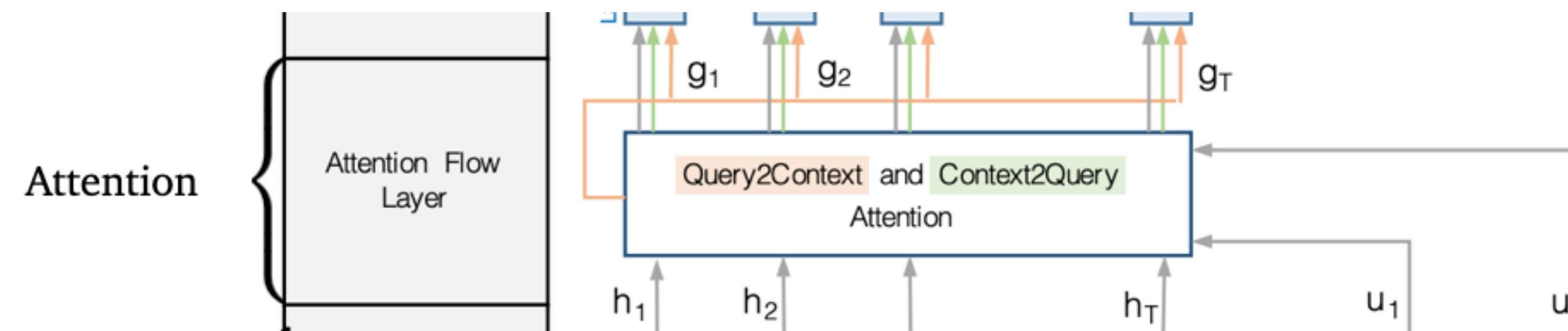
$$\mathbf{c}_i = [\vec{\mathbf{c}}_i; \overleftarrow{\mathbf{c}}_i] \in \mathbb{R}^{2H}$$

$$\vec{\mathbf{q}}_i = \text{LSTM}(\vec{\mathbf{q}}_{i-1}, e(q_i)) \in \mathbb{R}^H$$

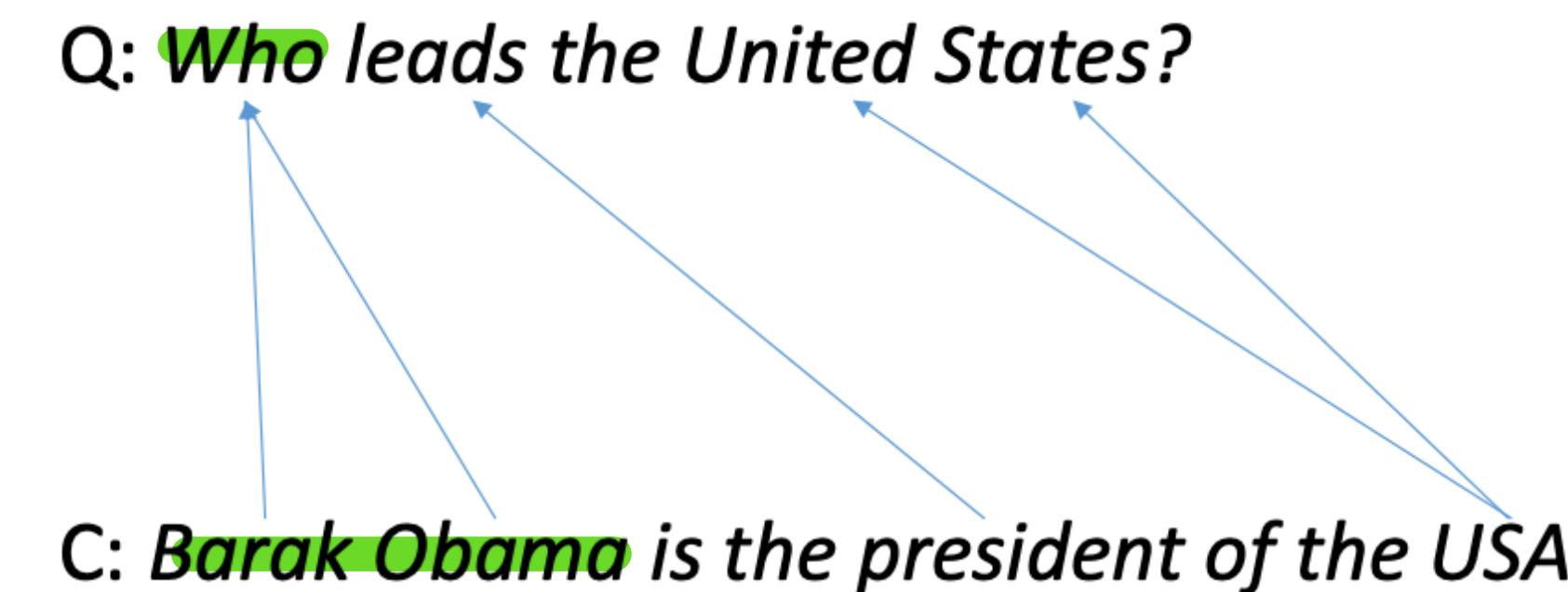
$$\overleftarrow{\mathbf{q}}_i = \text{LSTM}(\overleftarrow{\mathbf{q}}_{i+1}, e(q_i)) \in \mathbb{R}^H$$

$$\mathbf{q}_i = [\vec{\mathbf{q}}_i; \overleftarrow{\mathbf{q}}_i] \in \mathbb{R}^{2H}$$

# BiDAF: Attention

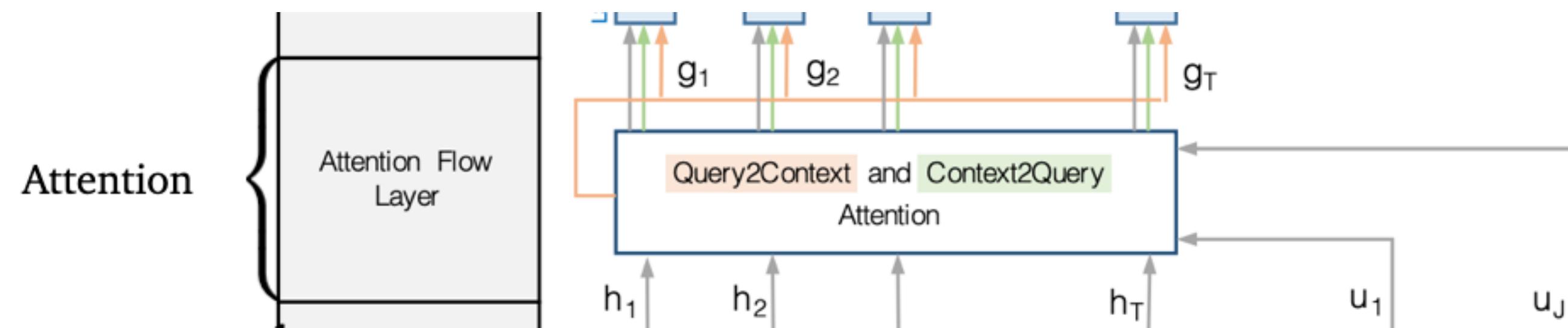


- Context-to-query attention: For each context word, choose the most relevant words from the query words.



For each context word, find the most relevant query word.

# BiDAF: Attention



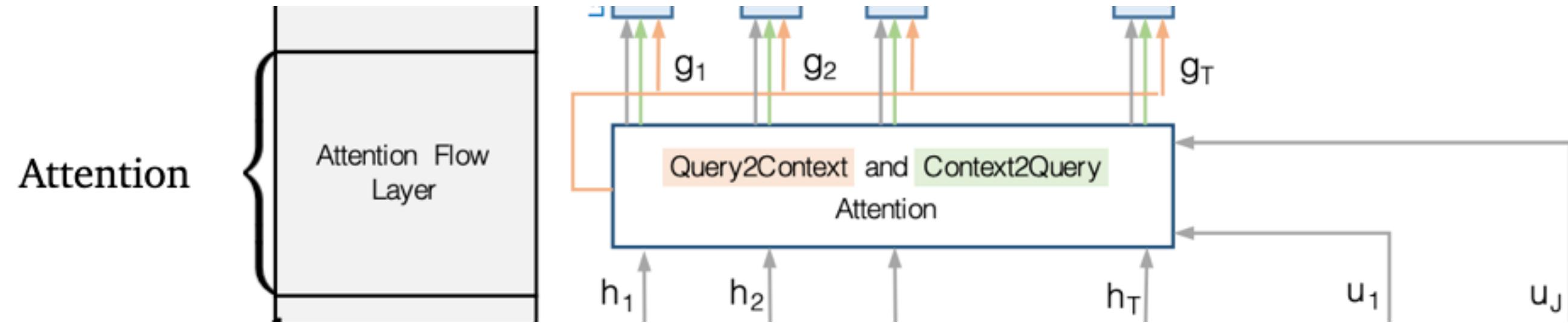
- **Query-to-context attention:** choose the context words that are most relevant to one of query words.

*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?*

# BiDAF: Attention

Key point is that we can try many things and keep what works and try to explain it



The final output is

$$\mathbf{g}_i = [\mathbf{c}_i; \mathbf{a}_i; \mathbf{c}_i \odot \mathbf{a}_i; \mathbf{c}_i \odot \mathbf{b}] \in \mathbb{R}^{8H}$$

- First, compute a similarity score for every pair of  $(\mathbf{c}_i, \mathbf{q}_j)$ :

$$S_{i,j} = \mathbf{w}_{\text{sim}}^T [\mathbf{c}_i; \mathbf{q}_j; \mathbf{c}_i \odot \mathbf{q}_j] \in \mathbb{R}^{2H, 2H, 2H}$$

$\mathbf{w}$  is a learnable parameter  
 $\mathbf{w}_{\text{sim}} \in \mathbb{R}^{6H}$

- Context-to-query attention (which question words are more relevant to  $c_i$ ):

$$\alpha_{i,j} = \text{softmax}_j(S_{i,j}) \in \mathbb{R}$$

$$\mathbf{a}_i = \sum_{j=1}^M \alpha_{i,j} \mathbf{q}_j \in \mathbb{R}^{2H}$$

weighted sum

- Query-to-context attention (which context words are relevant to some question words):

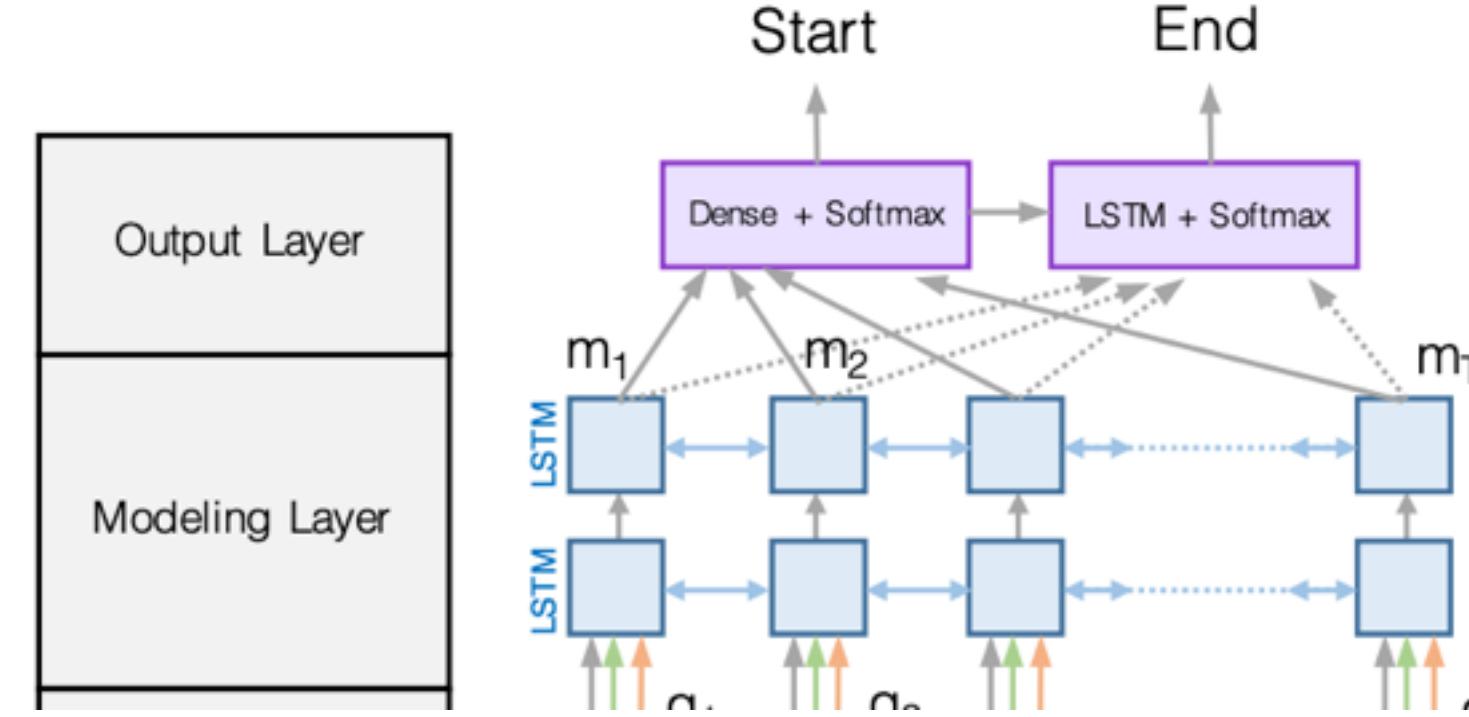
column vector with just context words in each row maybe helps n multiplication

$$\beta_i = \text{softmax}_i(\max_{j=1}^M (S_{i,j})) \in \mathbb{R}^N$$

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

a lot of words in context with low attention, just take max

# BiDAF: Modeling and output layers



The final training loss is

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

$-\log(p_{\text{start}} * p_{\text{end}})$

**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

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

**Output layer:** two classifiers predicting the start and end positions:

$$p_{\text{start}} = \text{softmax}(\mathbf{w}_{\text{start}}^\top [\mathbf{g}_i; \mathbf{m}_i]) \quad p_{\text{end}} = \text{softmax}(\mathbf{w}_{\text{end}}^\top [\mathbf{g}_i; \mathbf{m}'_i])$$

$$\mathbf{m}'_i = \text{BiLSTM}(\mathbf{m}_i) \in \mathbb{R}^{2H} \quad \mathbf{w}_{\text{start}}, \mathbf{w}_{\text{end}} \in \mathbb{R}^{10H}$$

# BiDAF: Performance on SQuAD

This model achieved 77.3 F1 on SQuAD v1.1.

- Without context-to-query attention  $\Rightarrow$  67.7 F1
- Without query-to-context attention  $\Rightarrow$  73.7 F1
- Without character embeddings  $\Rightarrow$  75.4 F1

	Published <sup>12</sup>	LeaderBoard <sup>13</sup>
Single Model	EM / F1	EM / F1
LR Baseline (Rajpurkar et al., 2016)	40.4 / 51.0	40.4 / 51.0
Dynamic Chunk Reader (Yu et al., 2016)	62.5 / 71.0	62.5 / 71.0
Match-LSTM with Ans-Ptr (Wang & Jiang, 2016)	64.7 / 73.7	64.7 / 73.7
Multi-Perspective Matching (Wang et al., 2016)	65.5 / 75.1	70.4 / 78.8
Dynamic Coattention Networks (Xiong et al., 2016)	66.2 / 75.9	66.2 / 75.9
FastQA (Weissenborn et al., 2017)	68.4 / 77.1	68.4 / 77.1
BiDAF (Seo et al., 2016)	68.0 / 77.3	68.0 / 77.3
SEDT (Liu et al., 2017a)	68.1 / 77.5	68.5 / 78.0
RaSoR (Lee et al., 2016)	70.8 / 78.7	69.6 / 77.7
FastQAExt (Weissenborn et al., 2017)	70.8 / 78.9	70.8 / 78.9
ReasoNet (Shen et al., 2017b)	69.1 / 78.9	70.6 / 79.4
Document Reader (Chen et al., 2017)	70.0 / 79.0	70.7 / 79.4
Ruminating Reader (Gong & Bowman, 2017)	70.6 / 79.5	70.6 / 79.5
jNet (Zhang et al., 2017)	70.6 / 79.8	70.6 / 79.8
Conductor-net	N/A	72.6 / 81.4
Interactive AoA Reader (Cui et al., 2017)	N/A	73.6 / 81.9
Reg-RaSoR	N/A	75.8 / 83.3
DCN+	N/A	74.9 / 82.8
AIR-FusionNet	N/A	76.0 / 83.9
R-Net (Wang et al., 2017)	72.3 / 80.7	76.5 / 84.3
BiDAF + Self Attention + ELMo	N/A	<b>77.9 / 85.3</b>
Reinforced Mnemonic Reader (Hu et al., 2017)	73.2 / 81.8	73.2 / 81.8

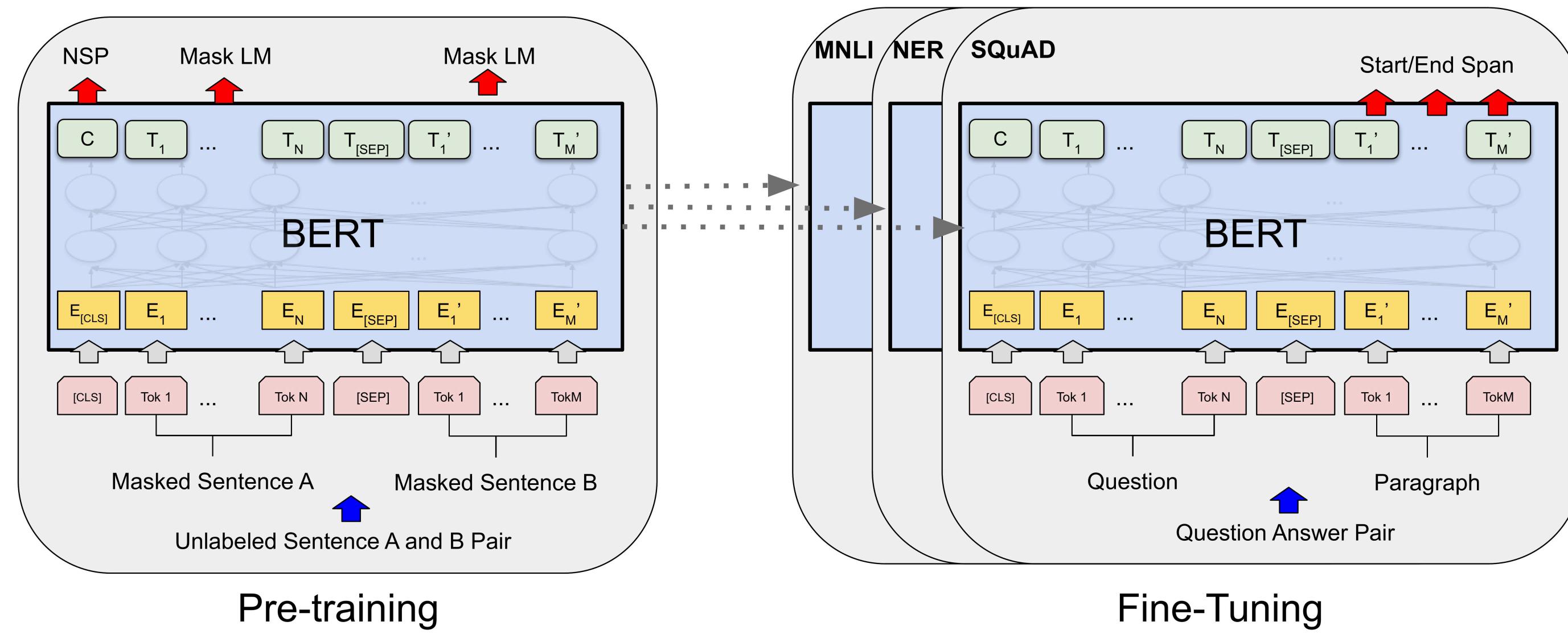
# Attention visualization

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 . As this was the 50th Super Bowl , the league emphasized the " golden anniversary " with various gold-themed initiatives , as well as temporarily suspending the tradition of naming each Super Bowl game with Roman numerals ( under which the game would have been known as " Super Bowl L " ) , so that the logo could prominently feature the Arabic numerals 50 .

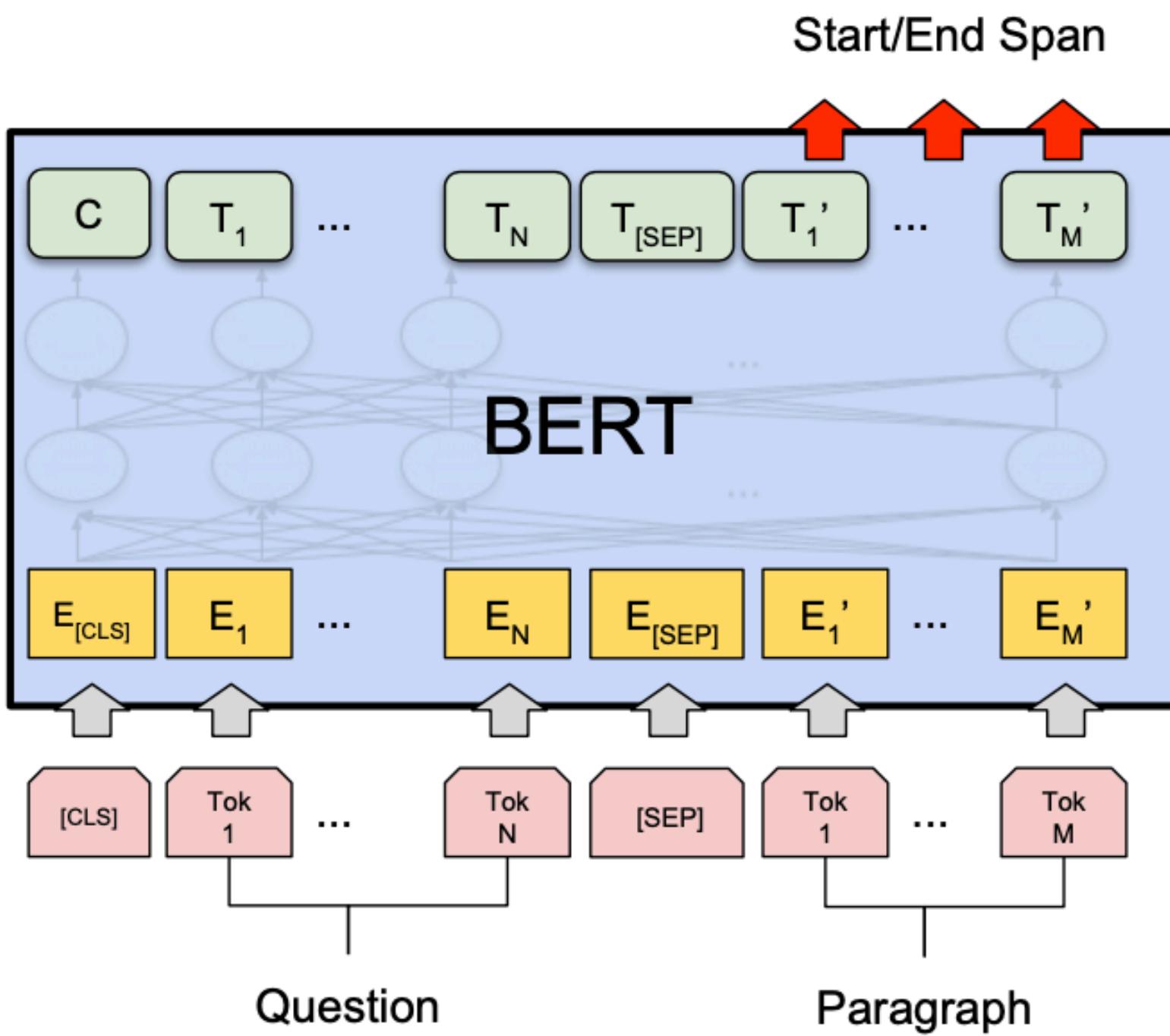


# BERT for reading comprehension

- BERT is a deep bidirectional Transformer encoder pre-trained on large amounts of text (Wikipedia + BooksCorpus)
- BERT is pre-trained on two training objectives:
  - Masked language model (MLM)
  - Next sentence prediction (NSP)
- $\text{BERT}_{\text{base}}$  has 12 layers and 110M parameters,  $\text{BERT}_{\text{large}}$  has 24 layers and 330M parameters



# BERT for reading comprehension



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

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

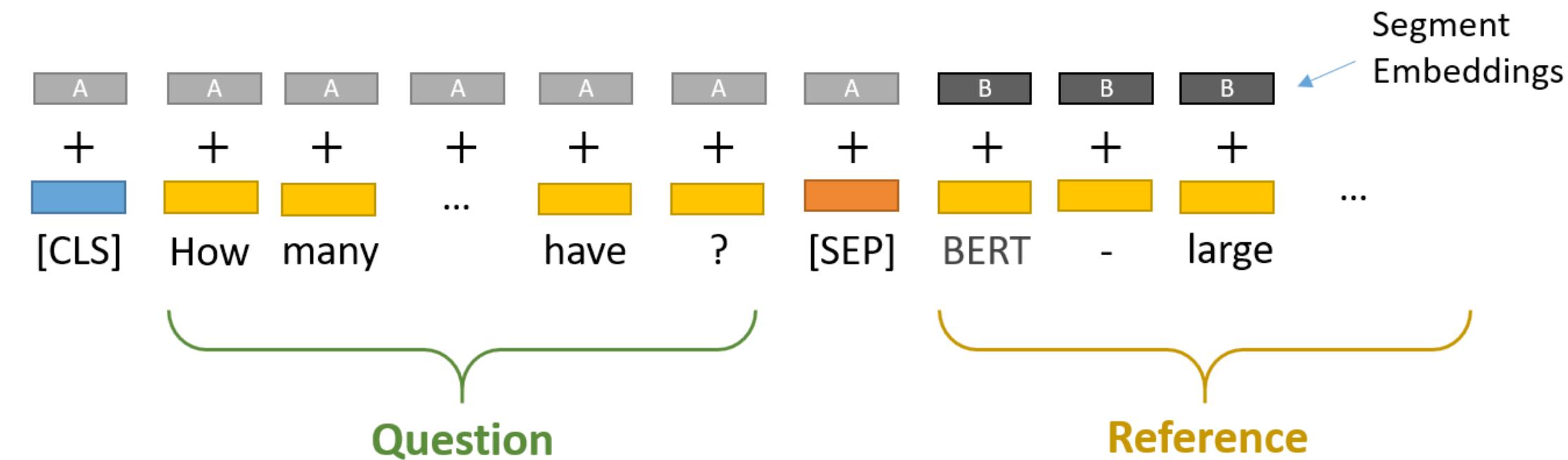
$$p_{\text{end}}(i) = \text{softmax}_i(\mathbf{w}_{\text{end}}^\top \mathbf{H})$$

where  $\mathbf{H} = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_N]$  are the hidden vectors of the paragraph, returned by BERT

**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/>

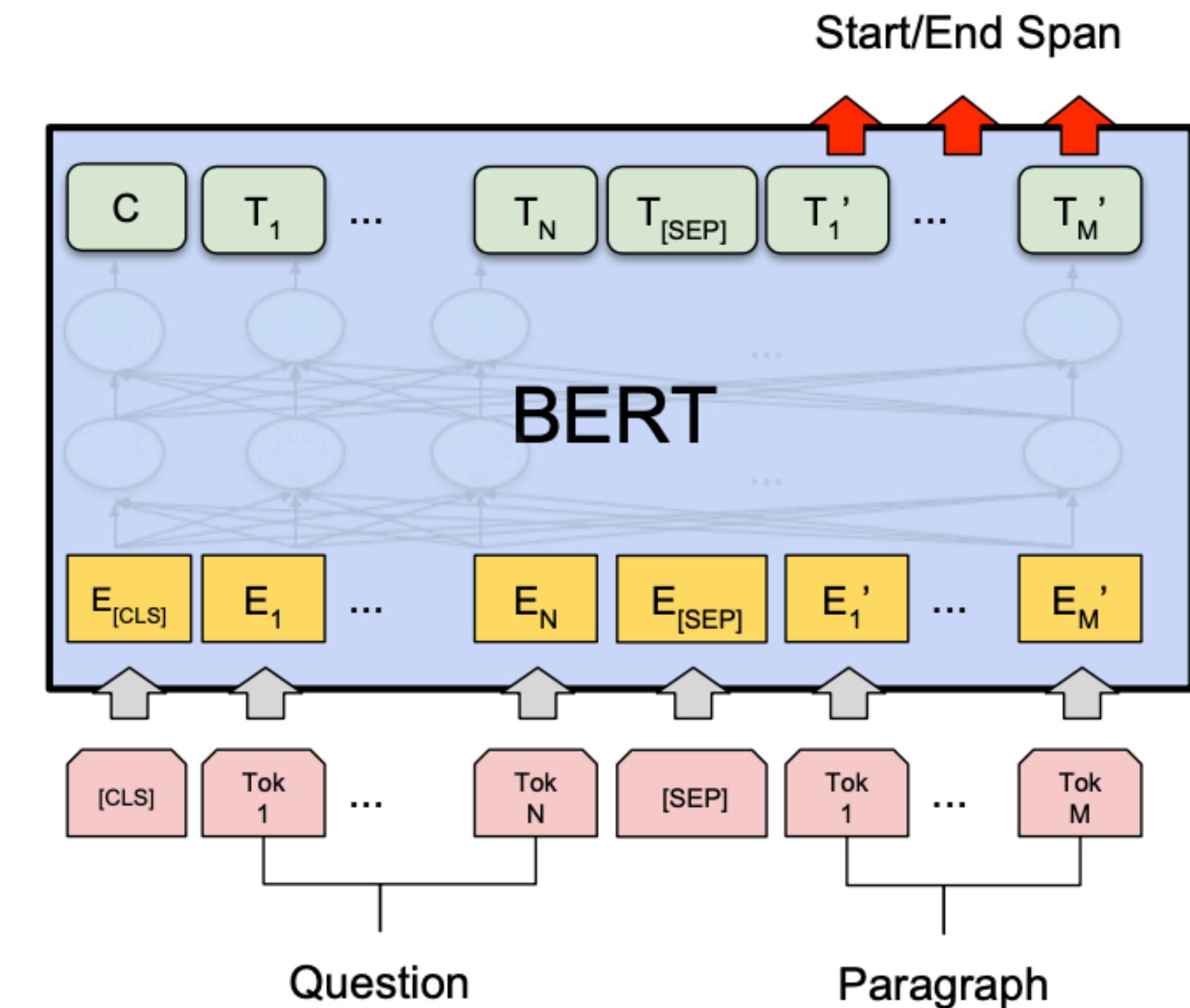
# BERT for reading comprehension

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

- All the BERT parameters (e.g., 110M) as well as the newly introduced parameters  $h_{\text{start}}, h_{\text{end}}$  (e.g.,  $768 \times 2 = 1536$ ) are optimized together for  $\mathcal{L}$ .
- It works amazingly well. Stronger pre-trained language models can lead to even better performance and SQuAD becomes a standard dataset for testing pre-trained models.

	F1	EM
Human performance	91.2*	82.3*
BiDAF	77.3	67.7
BERT-base	88.5	80.8
BERT-large	90.9	84.1
XLNet	94.5	89.0
RoBERTa	94.6	88.9
ALBERT	94.8	89.3

(dev set, except for human performance)



## Comparisons between BiDAF and BERT models

- BERT model has many many more parameters (110M or 330M) and BiDAF has ~2.5M parameters.
- BiDAF is built on top of several bidirectional LSTMs while BERT is built on top of Transformers (no recurrence architecture and easier to parallelize).
- BERT is **pre-trained** while BiDAF is only built on top of GloVe (and all the remaining parameters need to be learned from the supervision datasets).

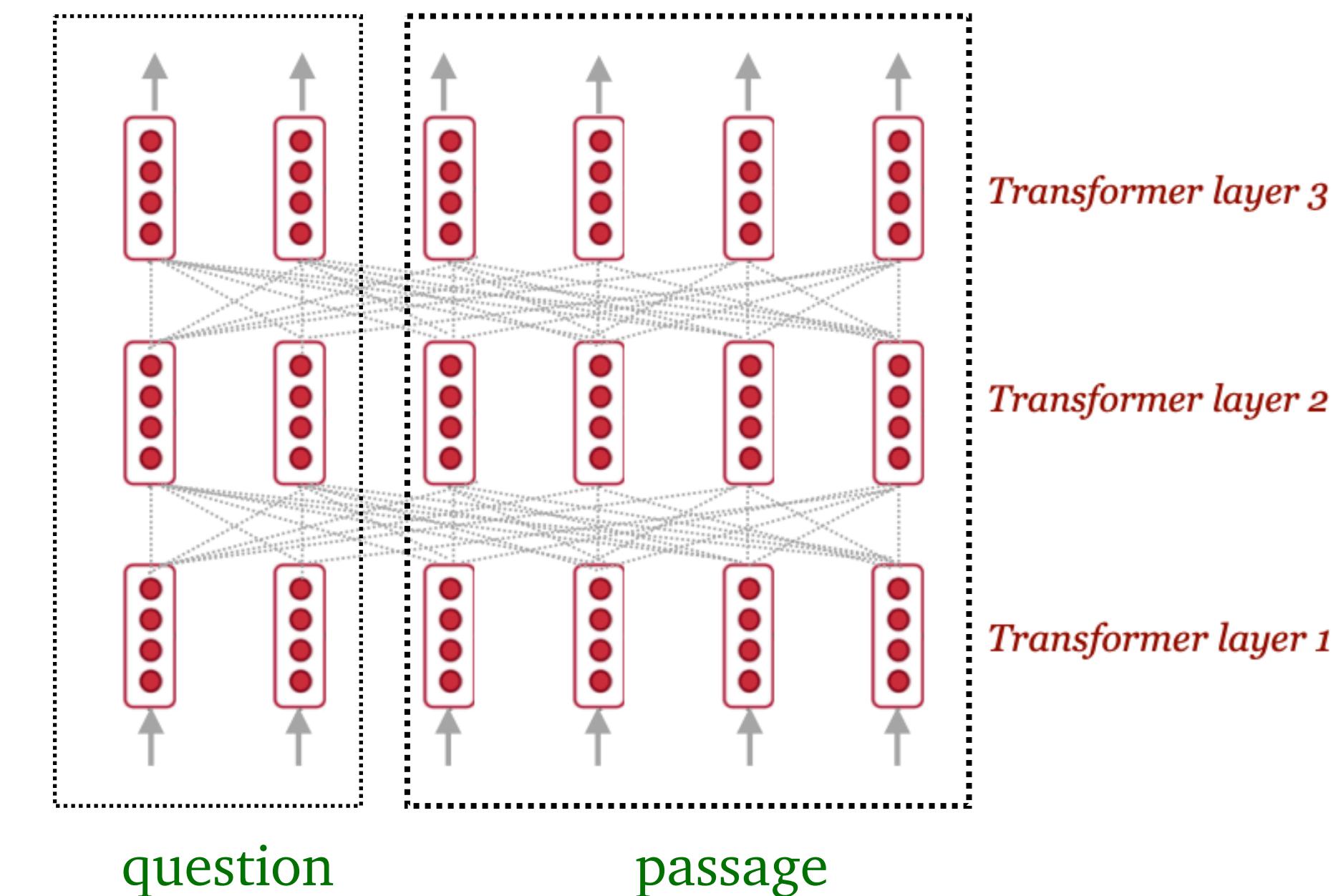
Pre-training is clearly a game changer but it is expensive..

# Comparisons between BiDAF and BERT models

Are they really fundamentally different? Probably not.

- BiDAF and other models aim to model the interactions between question and passage.
- BERT uses self-attention between the **concatenation** of question and passage =  
 $\text{attention}(P, P) + \text{attention}(P, Q) + \text{attention}(Q, P) + \text{attention}(Q, Q)$
- (Clark and Gardner, 2018) shows that adding a self-attention layer for the passage  
 $\text{attention}(P, P)$  to BiDAF also improves performance.

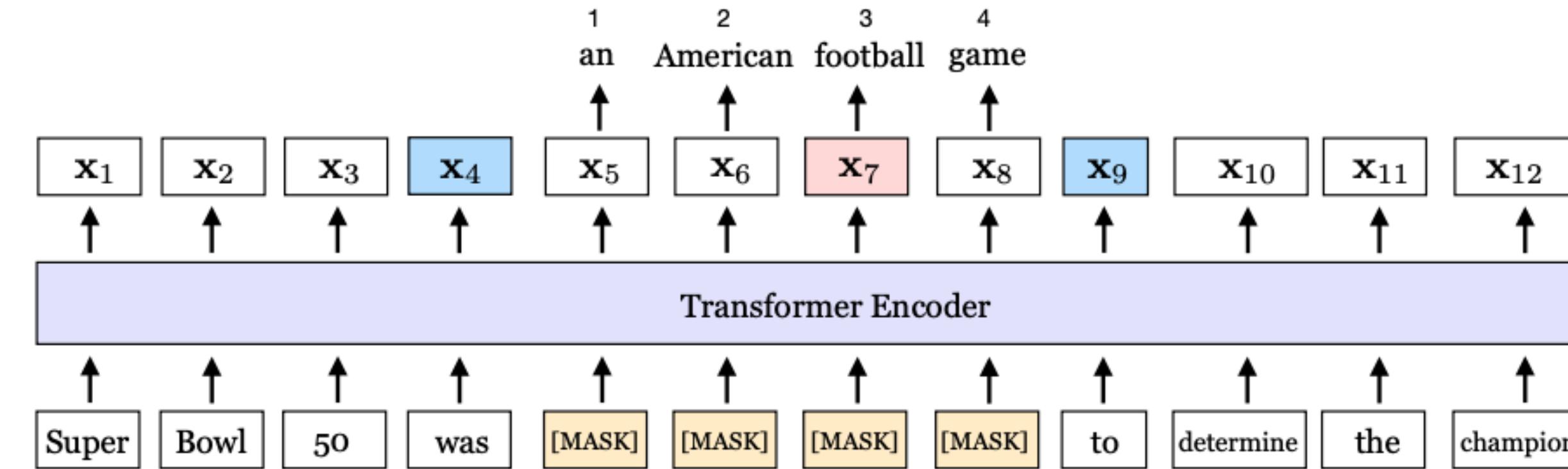
So what the BERT and BiDAF are doing is very similar



# Can we design better pre-training objectives?

The answer is yes!

$$\begin{aligned}\mathcal{L}(\text{football}) &= \mathcal{L}_{\text{MLM}}(\text{football}) + \mathcal{L}_{\text{SBO}}(\text{football}) \\ &= -\log P(\text{football} \mid \mathbf{x}_7) - \log P(\text{football} \mid \mathbf{x}_4, \mathbf{x}_9, \mathbf{p}_3)\end{aligned}$$



Two ideas:

- 1) masking contiguous spans of words instead of 15% random words
- 2) using the two end points of span to predict all the masked words in between = compressing the information of a span into its two endpoints

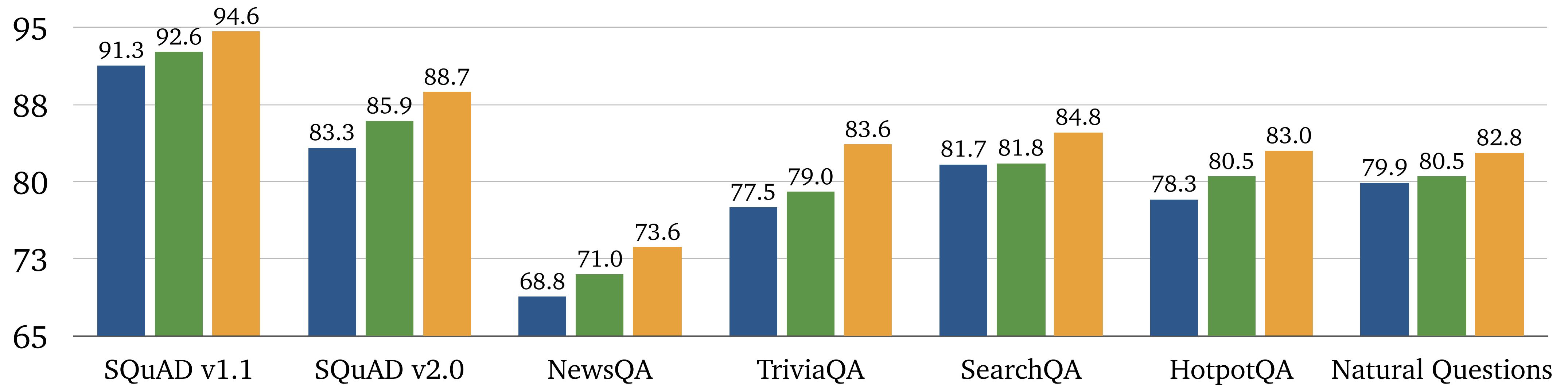
$$\begin{aligned}\mathbf{y}_i &= f(\mathbf{x}_{s-1}, \mathbf{x}_{e+1}, \mathbf{p}_{i-s+1}) \\ &\quad (\text{Endpoint, Endpoint, Positional Encoding}) \\ &\quad (\mathbf{x}_4, \mathbf{x}_9, \dots)\end{aligned}$$

# SpanBERT performance

Checkpoint of BERT from original paper(Google)

Google BERT  
Our BERT      BERT implemented by chen et.al.  
SpanBERT  
Spanbert

F1 scores



Without increasing more data or increasing model size, performance for question answering was improved by designing a better training objective

# Is reading comprehension solved?

- We have already surpassed human performance on SQuAD. Does it mean that reading comprehension is already solved? **Of course not!**
- The current systems still perform poorly on adversarial examples or examples from out-of-domain distributions

**Article:** Super Bowl 50

**Paragraph:** “*Peyton Manning became the first quarterback ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver’s Executive Vice President of Football Operations and General Manager. Quarterback Jeff Dean had jersey number 37 in Champ Bowl XXXIV.*”

**Question:** “What is the name of the quarterback who was 38 in Super Bowl XXXIII?”

**Original Prediction:** John Elway

**Prediction under adversary:** Jeff Dean

	Match Single	Match Ens.	BiDAF Single	BiDAF Ens.
Original	71.4	75.4	75.5	80.0
ADDSENT	27.3	29.4	34.3	34.2
ADDONESENT	39.0	41.8	45.7	46.9
ADDANY	7.6	11.7	4.8	2.7
ADDCOMMON	38.9	51.0	41.7	52.6

Adversarial example is an input designed to fool a machine learning model

# Is reading comprehension solved?

Systems trained on one dataset can't generalize to other datasets:

Fine-tuned on	Evaluated on				
	SQuAD	TriviaQA	NQ	QuAC	NewsQA
SQuAD	<b>75.6</b>	46.7	48.7	20.2	41.1
TriviaQA	49.8	<b>58.7</b>	42.1	20.4	10.5
NQ	53.5	46.3	<b>73.5</b>	21.6	24.7
QuAC	39.4	33.1	33.8	<b>33.3</b>	13.8
NewsQA	52.1	38.4	41.7	20.4	<b>60.1</b>

# Is reading comprehension solved?

## BERT-large model trained on SQuAD

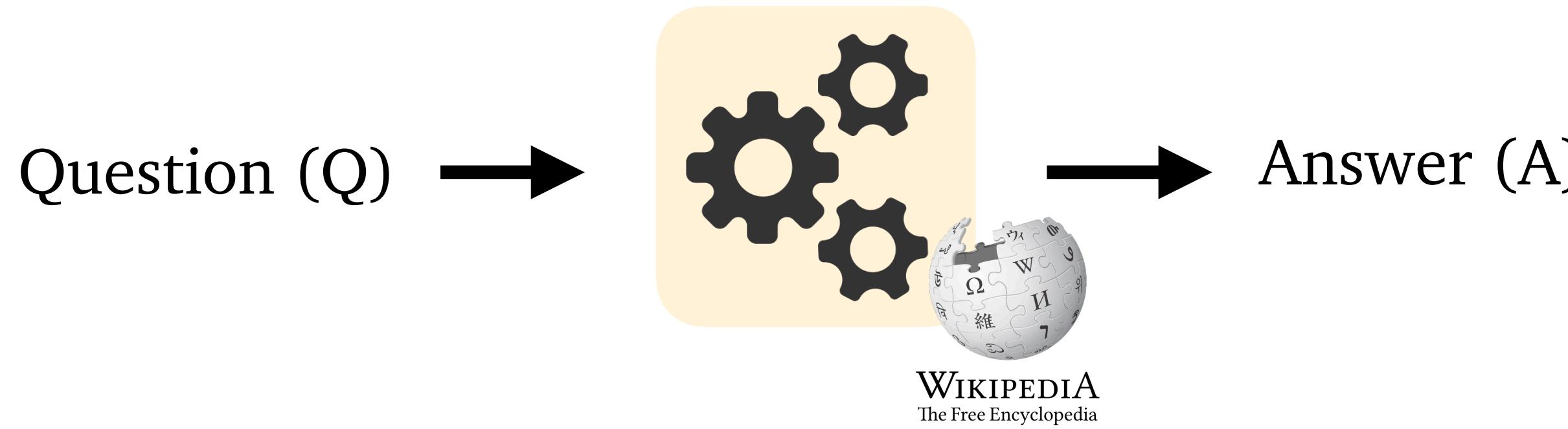
Test TYPE and Description	Failure Rate (%)	Example Test cases (with expected behavior and prediction)
Vocab	<i>MFT</i> : comparisons	20.0     C: Victoria is younger than Dylan. Q: Who is less young? A: Dylan ☑: Victoria
	<i>MFT</i> : intensifiers to superlative: most/least	91.3     C: Anna is worried about the project. Matthew is extremely worried about the project. Q: Who is least worried about the project? A: Anna ☑: Matthew
Taxonomy	<i>MFT</i> : match properties to categories	82.4     C: There is a tiny purple box in the room. Q: What size is the box? A: tiny ☑: purple
	<i>MFT</i> : nationality vs job	49.4     C: Stephanie is an Indian accountant. Q: What is Stephanie's job? A: accountant ☑: Indian accountant
	<i>MFT</i> : animal vs vehicles	26.2     C: Jonathan bought a truck. Isabella bought a hamster. Q: Who bought an animal? A: Isabella ☑: Jonathan
	<i>MFT</i> : comparison to antonym	67.3     C: Jacob is shorter than Kimberly. Q: Who is taller? A: Kimberly ☑: Jacob
	<i>MFT</i> : more/less in context, more/less antonym in question	100.0     C: Jeremy is more optimistic than Taylor. Q: Who is more pessimistic? A: Taylor ☑: Jeremy
Robust.	<i>INV</i> : Swap adjacent characters in Q (typo)	11.6     C: ...Newcomen designs had a duty of about 7 million, but most were closer to 5 million.... Q: What was the ideal duty → udy of a Newcomen engine? A: INV ☑: 7 million → 5 million
	<i>INV</i> : add irrelevant sentence to C	9.8     (no example)

# Is reading comprehension solved?

## BERT-large model trained on SQuAD

Temporal	<b>MFT:</b> change in one person only	41.5	C: Both Luke and Abigail were writers, but there was a change in Abigail, who is now a model. Q: Who is a model? A: Abigail  Abigail were writers, but there was a change in Abigail
	<b>MFT:</b> Understanding before/after, last/first	82.9	C: Logan became a farmer before Danielle did. Q: Who became a farmer last? A: Danielle  Logan
Neg.	<b>MFT:</b> Context has negation	67.5	C: Aaron is not a writer. Rebecca is. Q: Who is a writer? A: Rebecca  Aaron
	<b>MFT:</b> Q has negation, C does not	100.0	C: Aaron is an editor. Mark is an actor. Q: Who is not an actor? A: Aaron  Mark
Coref.	<b>MFT:</b> Simple coreference, he/she.	100.0	C: Melissa and Antonio are friends. He is a journalist, and she is an adviser. Q: Who is a journalist? A: Antonio  Melissa
	<b>MFT:</b> Simple coreference, his/her.	100.0	C: Victoria and Alex are friends. Her mom is an agent Q: Whose mom is an agent? A: Victoria  Alex
SRL	<b>MFT:</b> former/latter	100.0	C: Kimberly and Jennifer are friends. The former is a teacher Q: Who is a teacher? A: Kimberly  Jennifer
	<b>MFT:</b> subject/object distinction	60.8	C: Richard bothers Elizabeth. Q: Who is bothered? A: Elizabeth  Richard
	<b>MFT:</b> subj/obj distinction with 3 agents	95.7	C: Jose hates Lisa. Kevin is hated by Lisa. Q: Who hates Kevin? A: Lisa  Jose

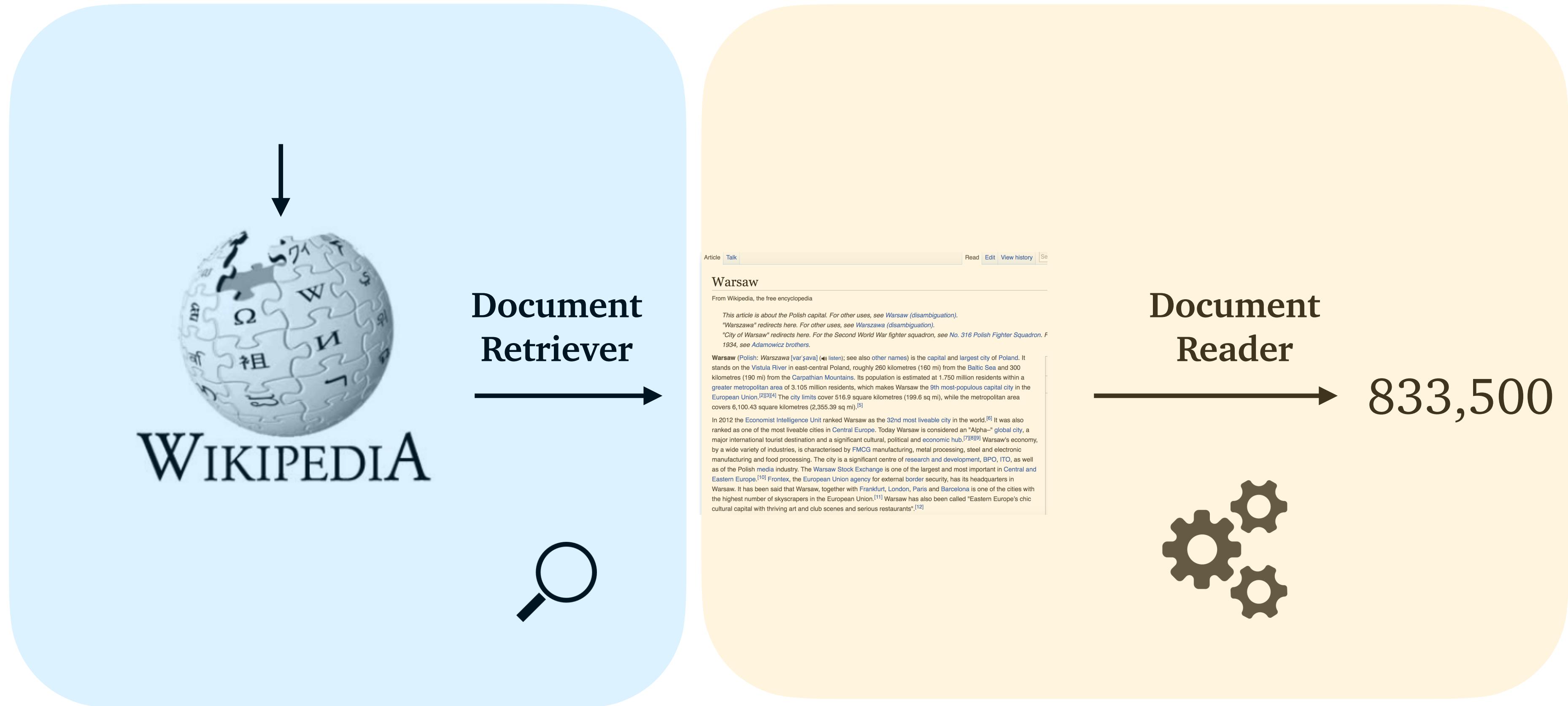
### 3. Open-domain question answering



- Different from reading comprehension, we don't assume a given passage.
- Instead, we only have access to a large collection of documents (e.g., Wikipedia). 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!

*In contrast to **closed-domain** systems that deal with questions under a specific domain (medicine, technical support)..*

# Retriever-reader framework



# Retriever-reader framework

- Input: a large collection of documents  $\mathcal{D} = D_1, D_2, \dots, D_N$  and  $Q$
  - Output: an answer string  $A$
- 
- Retriever:  $f(\mathcal{D}, Q) \rightarrow P_1, \dots, P_K$       K is pre-defined (e.g., 100)
  - Reader:     $g(Q, \{P_1, \dots, P_K\}) \rightarrow A$       A reading comprehension problem!

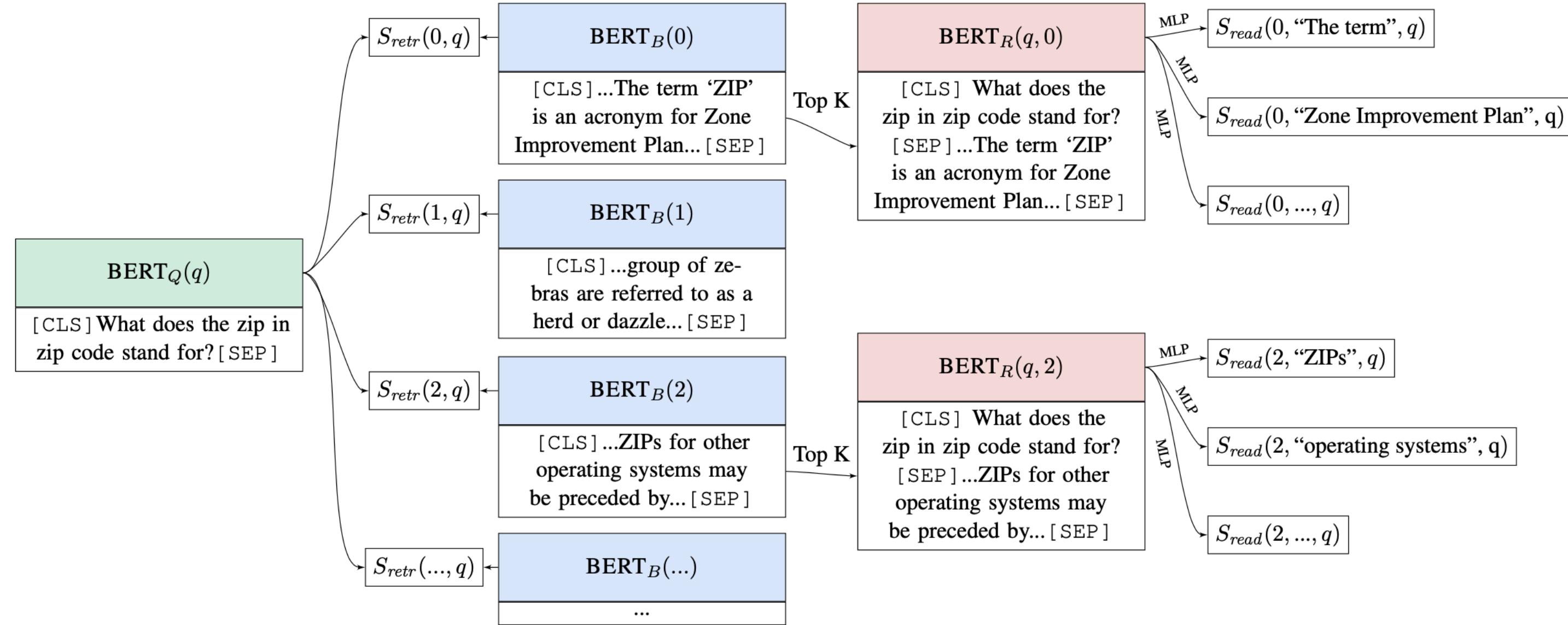
In DrQA,

- Retriever = A standard TF-IDF information-retrieval sparse model (a fixed module)
- Reader = a neural reading comprehension model that we just learned
  - Trained on SQuAD and other distantly-supervised QA datasets

*Distantly-supervised examples:*  $(Q, A) \rightarrow (P, Q, A)$

# We can train the retriever too

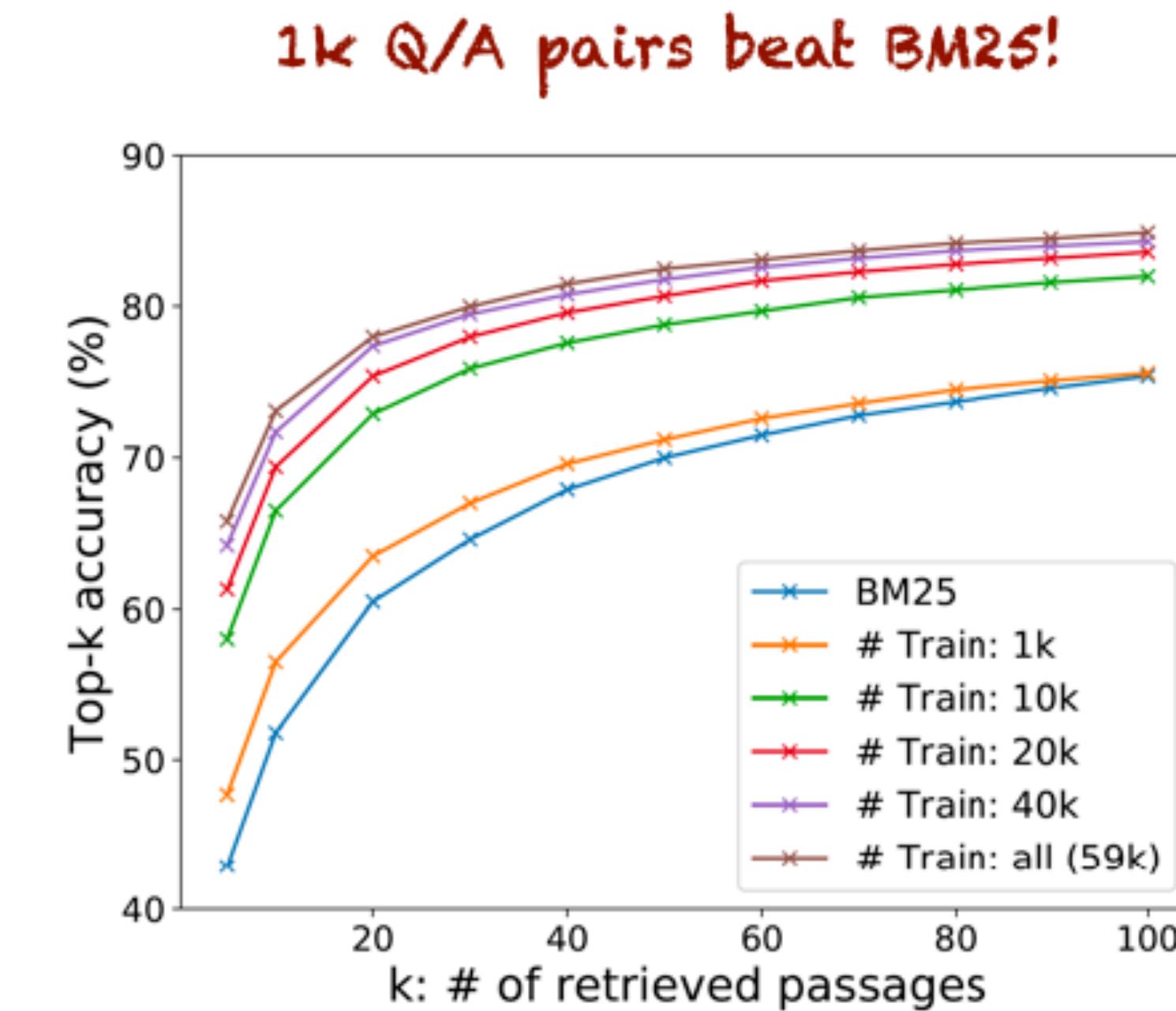
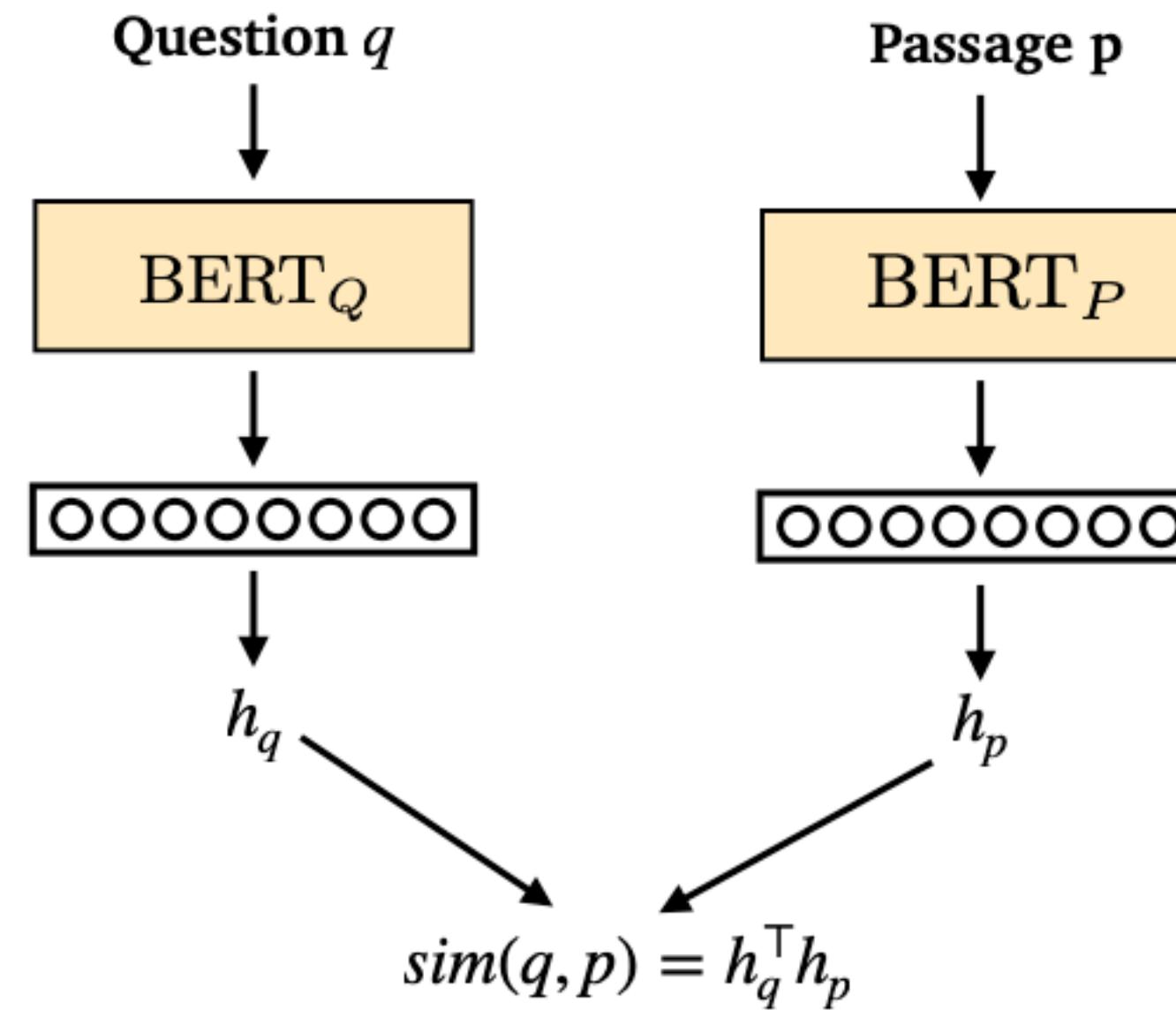
- Joint training of retriever and reader



- Each text passage can be encoded as a vector using BERT and the retriever score can be measured as the dot product between the question representation and passage representation.
- However, it is not easy to model as there are a huge number of passages (e.g., 21M in English Wikipedia)

# We can train the retriever too

- Dense passage retrieval (DPR) - We can also just train the retriever using question-answer pairs!



- Trainable retriever (using BERT) largely outperforms traditional IR retrieval models

# We can train the retriever too

Who tells harry potter that he is a wizard in the harry potter series? ▼ Run

Title: *Harry Potter (film series)*      Retrieval ranking: #90       $P(p|q)=0.85$      $P(a|p,q)=1.00$      $P(a,p|q)=0.84$

... and uncle. At the age of eleven, half-giant **Rubeus Hagrid** informs him that he is actually a wizard and that his parents were murdered by an evil wizard named Lord Voldemort. Voldemort also attempted to kill one-year-old Harry on the same night, but his killing curse mysteriously rebounded and reduced him to a weak and helpless form. Harry became extremely famous in the Wizarding World as a result. Harry begins his first year at Hogwarts School of Witchcraft and Wizardry and learns about magic. During the year, Harry and his friends Ron Weasley and Hermione Granger become entangled in the ...

Title: *Harry Potter (character)*      Retrieval ranking: #1       $P(p|q)=0.04$      $P(a|p,q)=0.97$      $P(a,p|q)=0.04$

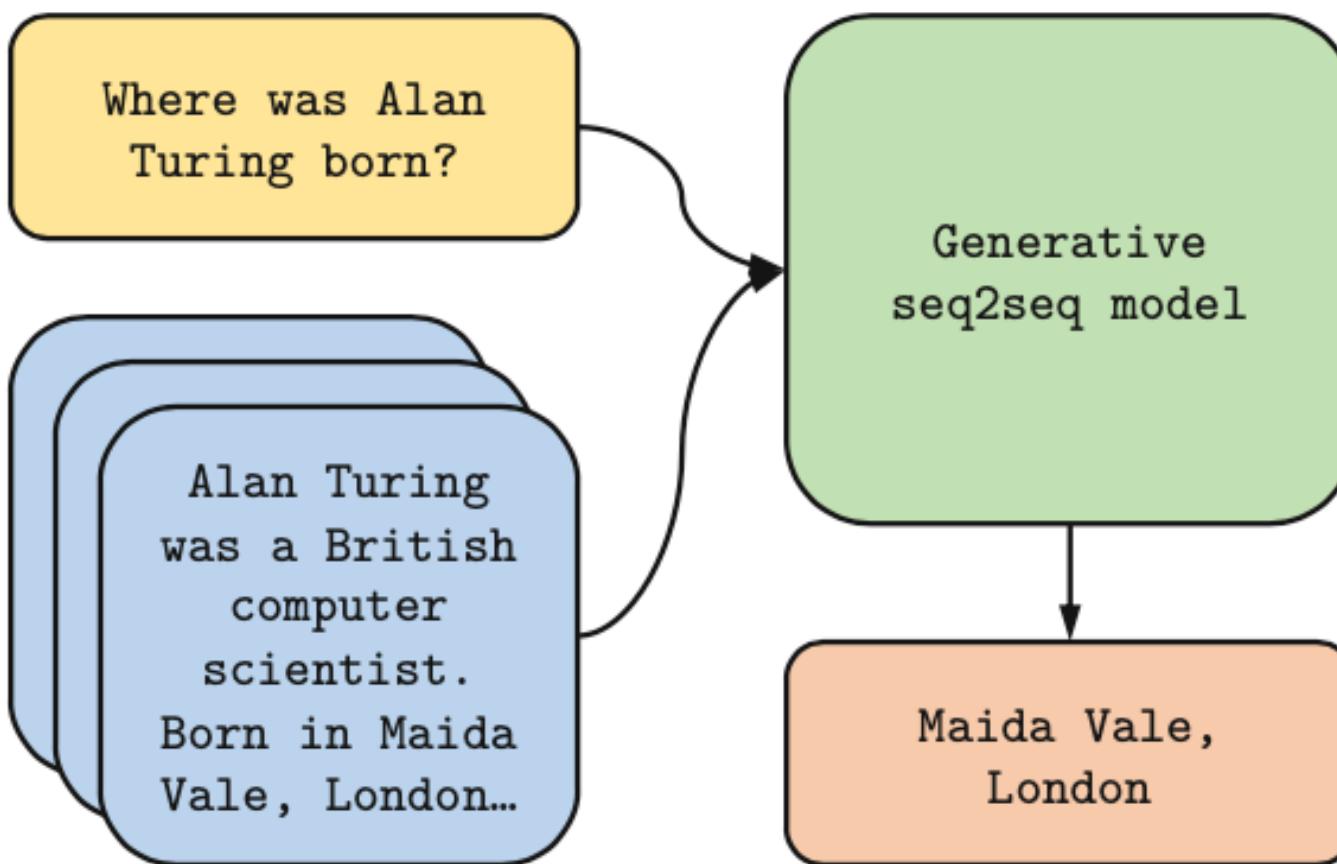
... Harry Potter (character) Harry James Potter is the titular protagonist of J. K. Rowling's "Harry Potter" series. The majority of the books' plot covers seven years in the life of the orphan Potter, who, on his eleventh birthday, learns he is a wizard. Thus, he attends Hogwarts School of Witchcraft and Wizardry to practice magic under the guidance of the kindly headmaster Albus Dumbledore and other school professors along with his best friends Ron Weasley and **Hermione Granger**. Harry also discovers that he is already famous throughout the novel's magical community, and that his fate is tied with that of ...

<http://qa.cs.washington.edu:2020/>

# Dense retrieval + generative models

Recent work shows that it is beneficial to generate answers instead of to extract answers.

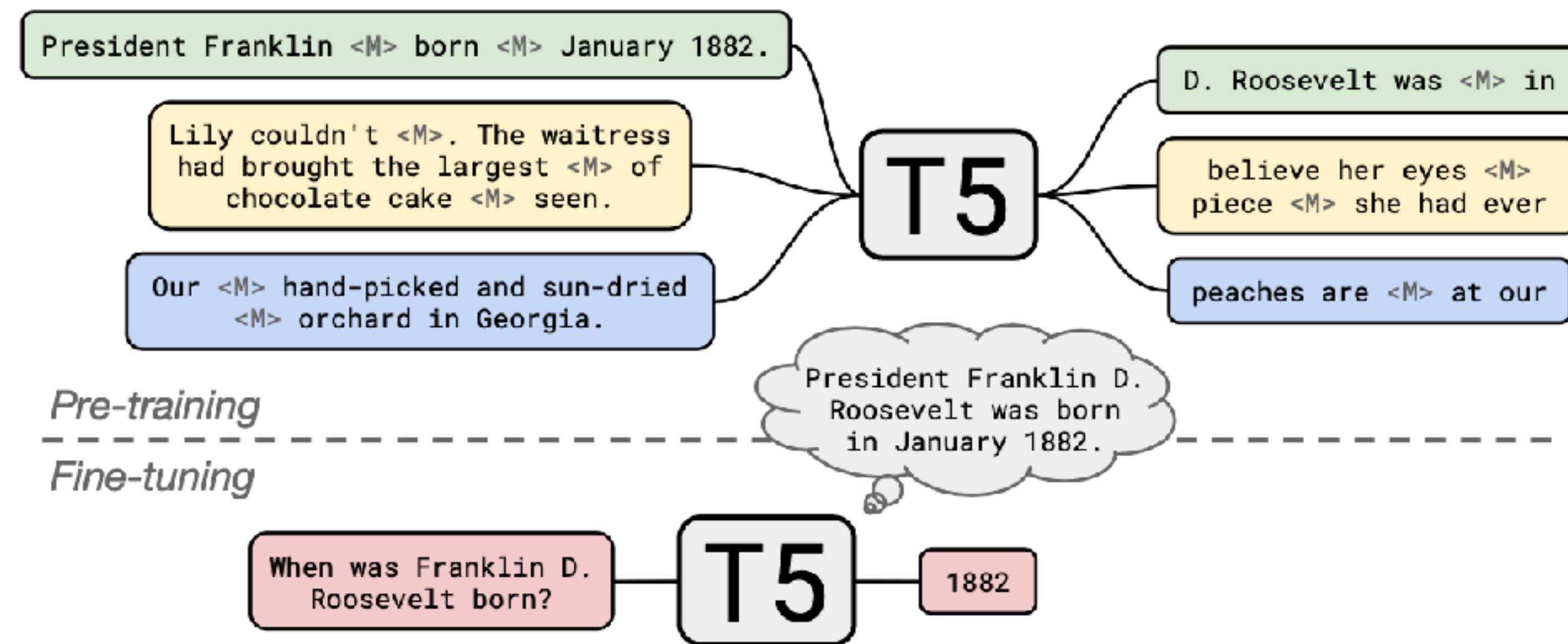
Fusion-in-decoder (FID) = DPR + T5



Model	NaturalQuestions	TriviaQA	
ORQA (Lee et al., 2019)	31.3	45.1	-
REALM (Guu et al., 2020)	38.2	-	-
DPR (Karpukhin et al., 2020)	41.5	57.9	-
SpanSeqGen (Min et al., 2020)	42.5	-	-
RAG (Lewis et al., 2020)	44.5	56.1	68.0
T5 (Roberts et al., 2020)	36.6	-	60.5
GPT-3 few shot (Brown et al., 2020)	29.9	-	71.2
Fusion-in-Decoder (base)	48.2	65.0	77.1
Fusion-in-Decoder (large)	<b>51.4</b>	<b>67.6</b>	<b>80.1</b>

# Large language models can do open-domain QA well

... without an explicit retriever stage



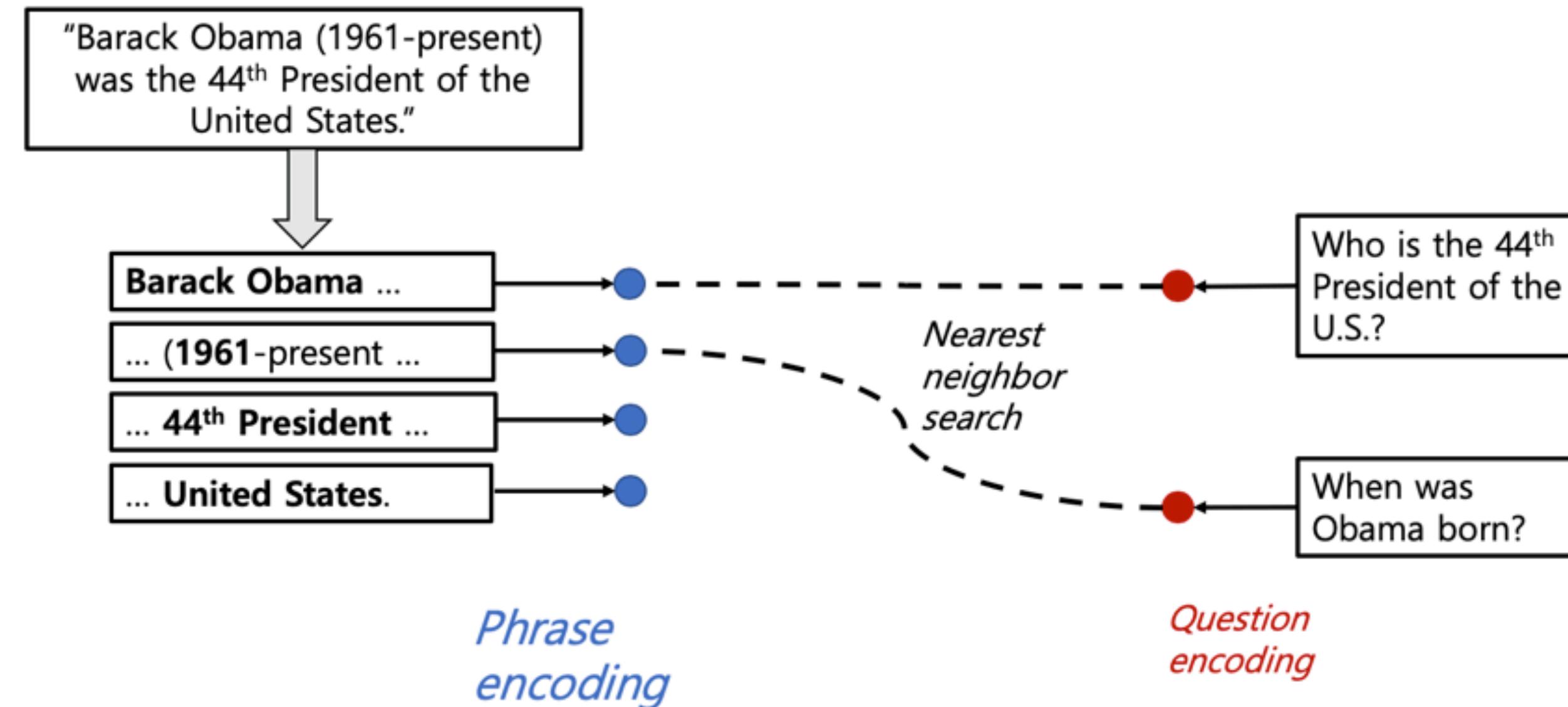
# Maybe the reader model is not necessary too!

Just retrieve the phrase no need to read

Faster than using bert models

It is possible to encode all the phrases (60 billion phrases in Wikipedia) using dense vectors and only do nearest neighbor search without a BERT model at inference time!

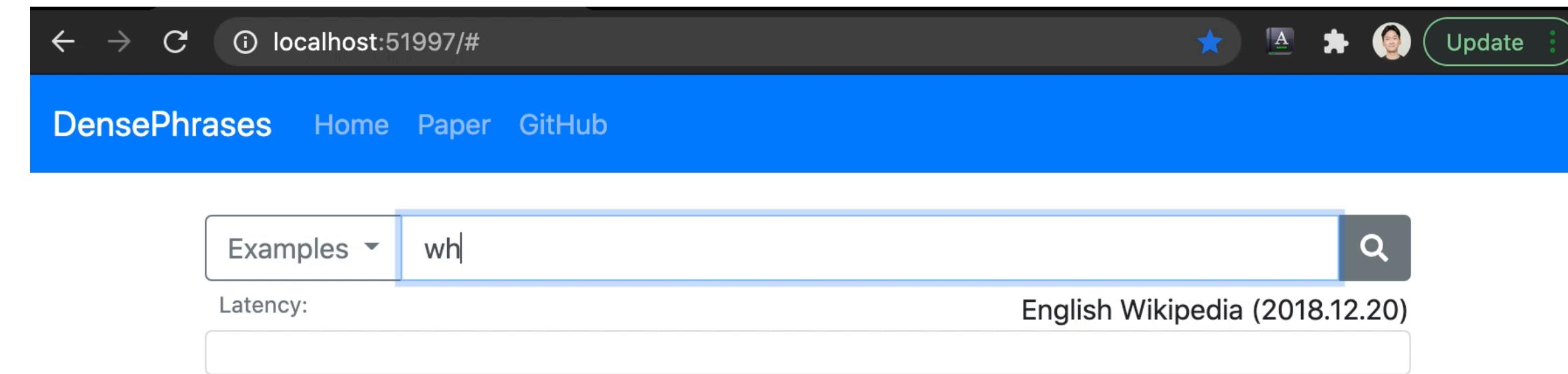
## Phrase Indexing



Seo et al., 2019. Real-Time Open-Domain Question Answering with Dense-Sparse Phrase Index

Lee et al., 2020. Learning Dense Representations of Phrases at Scale

# DensePhrases: Demo



# Thanks!

[danqic@cs.princeton.edu](mailto:danqic@cs.princeton.edu)