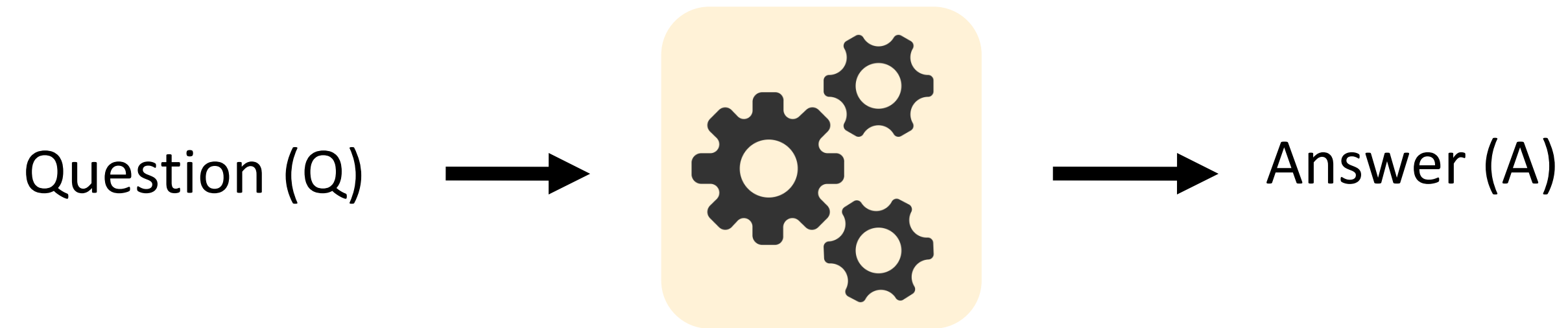# 1. What is question answering?
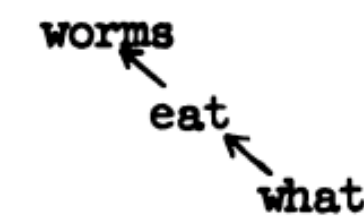
Question (Q) ➡️ ⚙️ ➡️ Answer (A)

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!
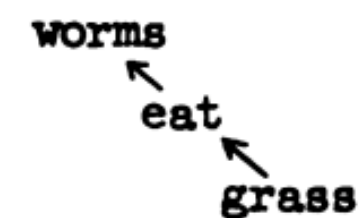
(Simmons et al., 1964)
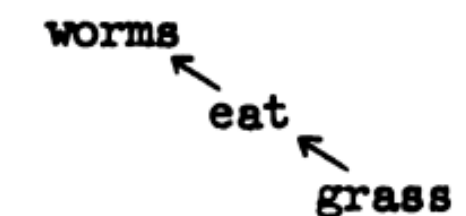
```
Question:
            a)  What do worms eat?
                    worms
                      ↖
                       eat
                         ↖
                          what
═══════════════════════════════════════════════
Answers:
 b)  Worms eat grass          c)  Grass is eaten by worms
        worms                       → worms eat grass
          ↖                              worms
           eat                             ↖
             ↖                              eat
              grass                           ↖
                                               grass
          (complete agreement of dependencies)
```
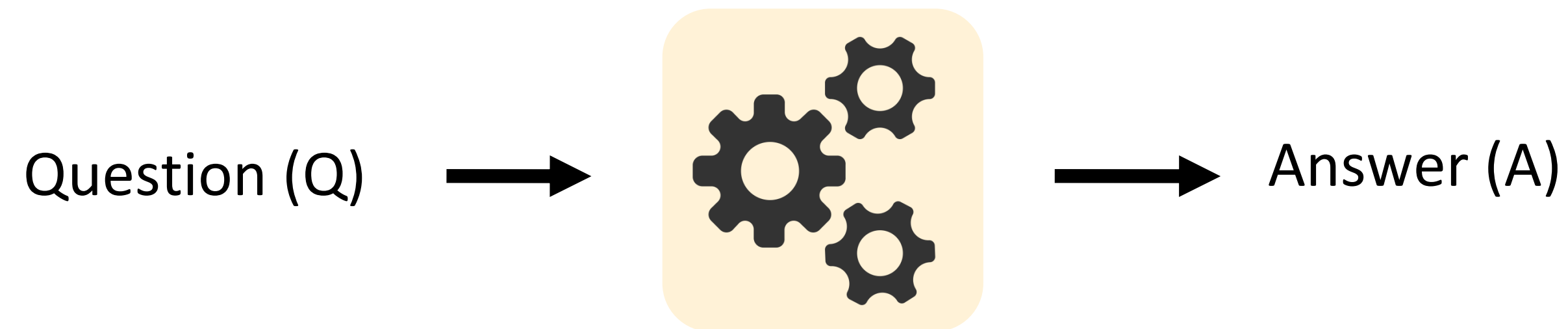
# Question answering: a taxonomy

Question (Q) ⟶  ⟶ Answer (A)

- 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, …

# 2. Reading comprehension

**Reading comprehension =** comprehend a passage of text and answer questions about its content  (P, Q) $\longrightarrow$ 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 =** comprehend a passage of text and answer questions about its content  (P, Q) $\longrightarrow$ 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 .

**finished**
- 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

**beating**
- 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 was for years the most popular reading comprehension dataset; it is "almost solved" today (though the underlying task is not,) 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**

(Rajpurkar et al., 2016): SQuAD: 100,000+ Questions for Machine Comprehension of Text

# 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

# Other question answering datasets

- TriviaQA: Questions and answers by trivia enthusiasts. Independently collected web paragraphs that contain the answer and seem to discuss question, but no human verification that paragraph supports answer to question

- Natural Questions: Question drawn from frequently asked Google search questions. Answers from Wikipedia paragraphs. Answer can be substring, yes, no, or NOT_PRESENT. Verified by human annotation.

- HotpotQA. Constructed questions to be answered from the whole of Wikipedia which involve getting information from two pages to answer a multistep query:
  Q: Which novel by the author of "Armada" will be adapted as a feature film by Steven Spielberg?  A: *Ready Player One*
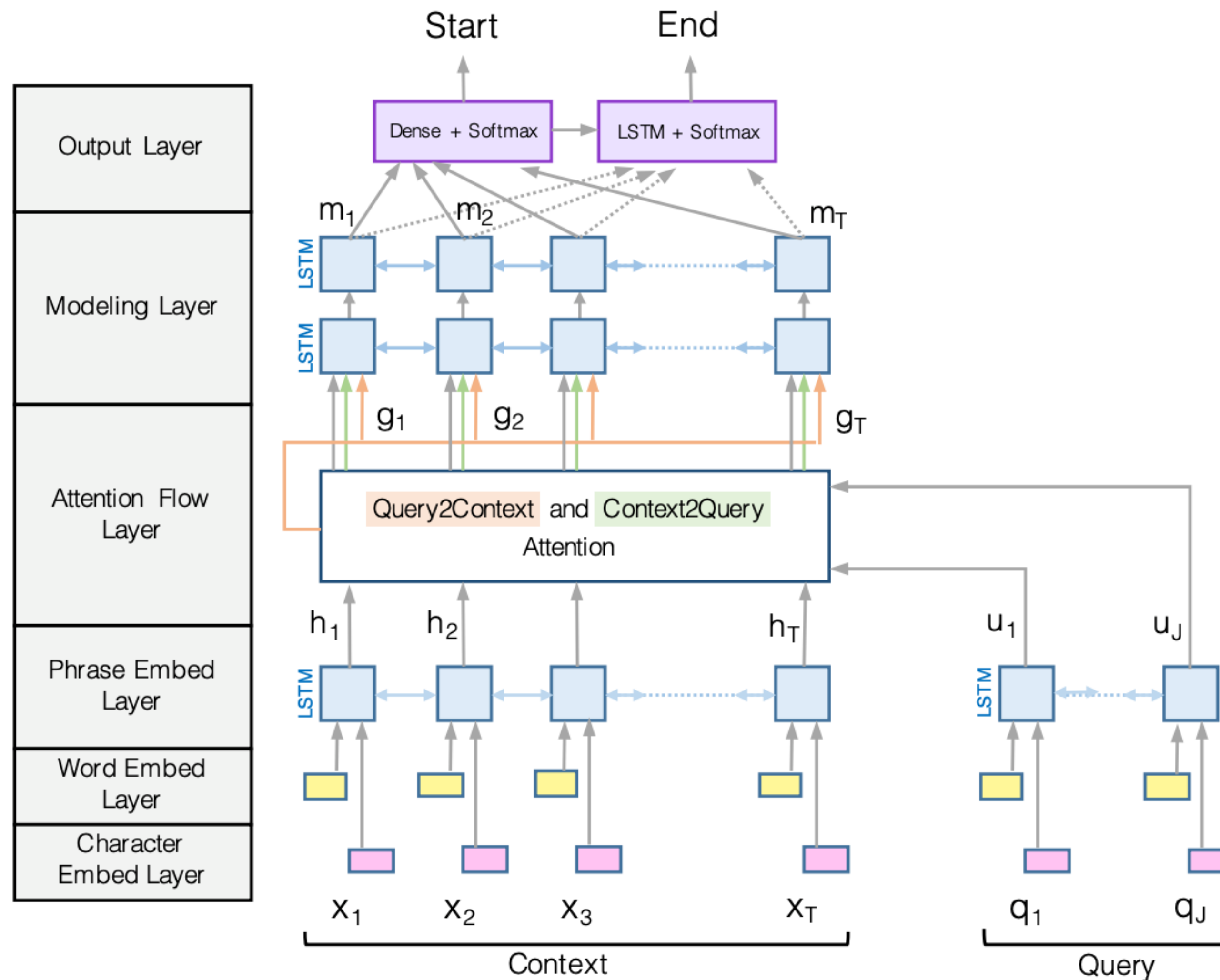
# Neural models for reading comprehension
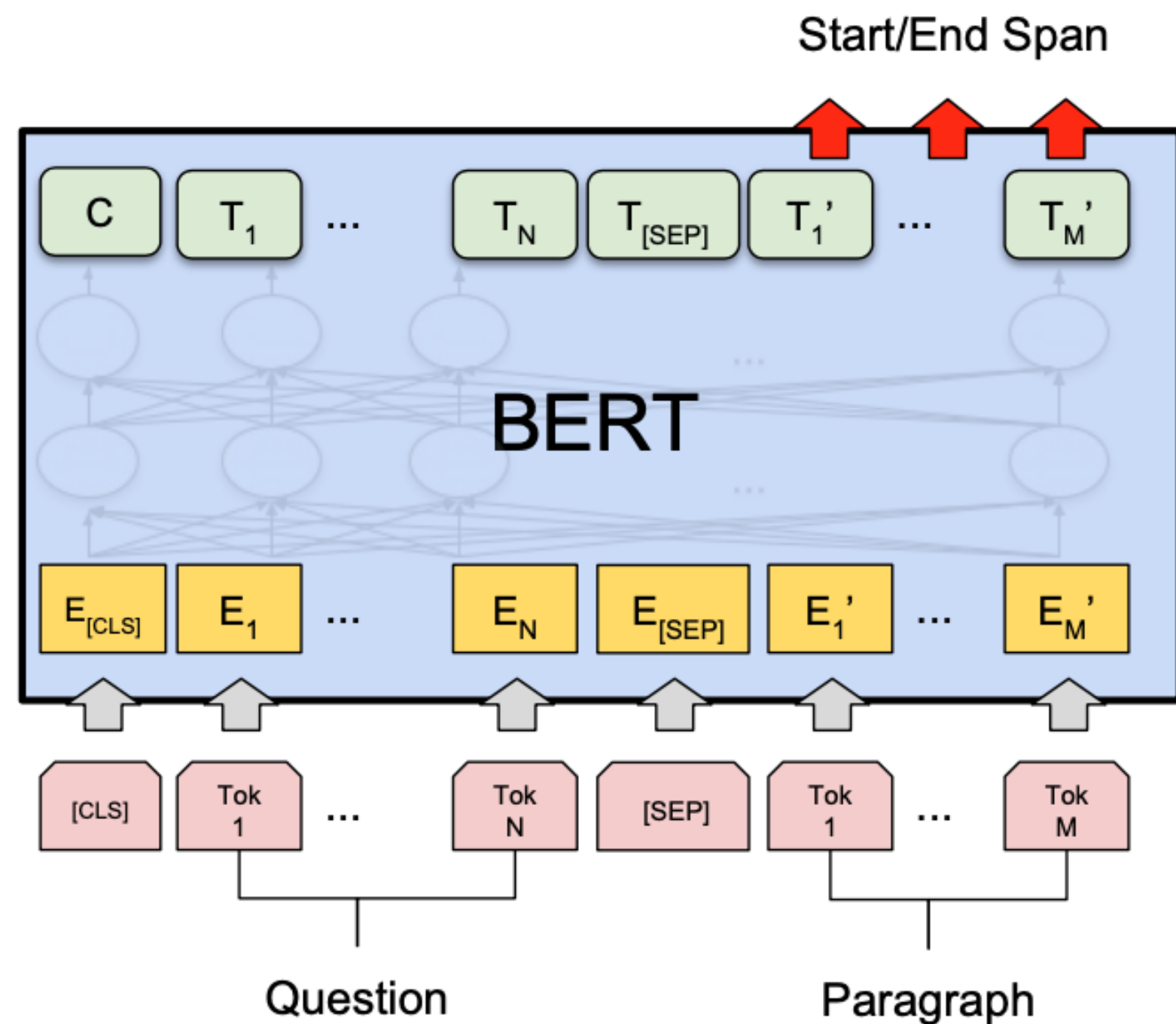
**How can we build a model to solve SQuAD?**

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

- Problem formulation
  - Input: $C = (c_1, c_2, \ldots, c_N), Q = (q_1, q_2, \ldots, q_M), c_i, q_i \in V$     N~100, M ~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), BiDAF (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+)
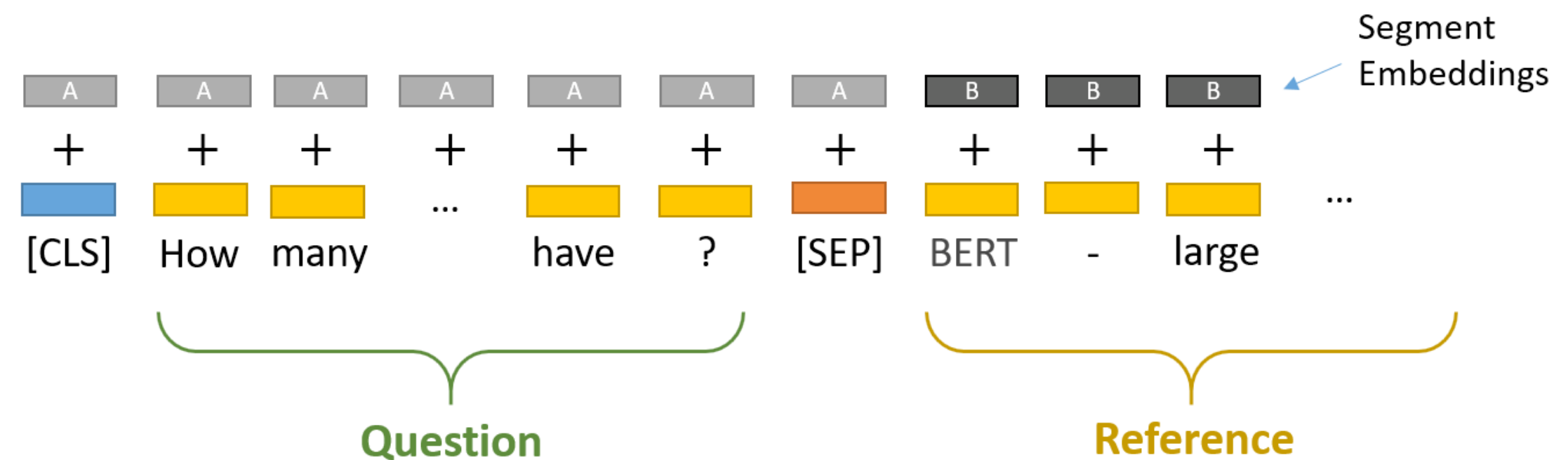
# BiDAF: the Bidirectional Attention Flow model

(Seo et al., 2017): Bidirectional Attention Flow for Machine Comprehension

# BERT for reading comprehension



Start/End Span

**BERT**

Question   Paragraph

**Question** = Segment A

**Passage** = Segment B

**Answer** = predicting two endpoints in segment B

Segment Embeddings

| A | A | A | A | A | A | A | B | B | B |
|---|---|---|---|---|---|---|---|---|---|

[CLS]  How  many  ...  have  ?  [SEP]  BERT  -  large  ...

**Question**          **Reference**

**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/

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

$$p_{\text{start}}(i) = \text{softmax}_i(\mathbf{w}_{\text{start}}^\mathsf{T} \mathbf{h}_i)$$

$$p_{\text{end}}(i) = \text{softmax}_i(\mathbf{w}_{\text{end}}^\mathsf{T} \mathbf{h}_i)$$

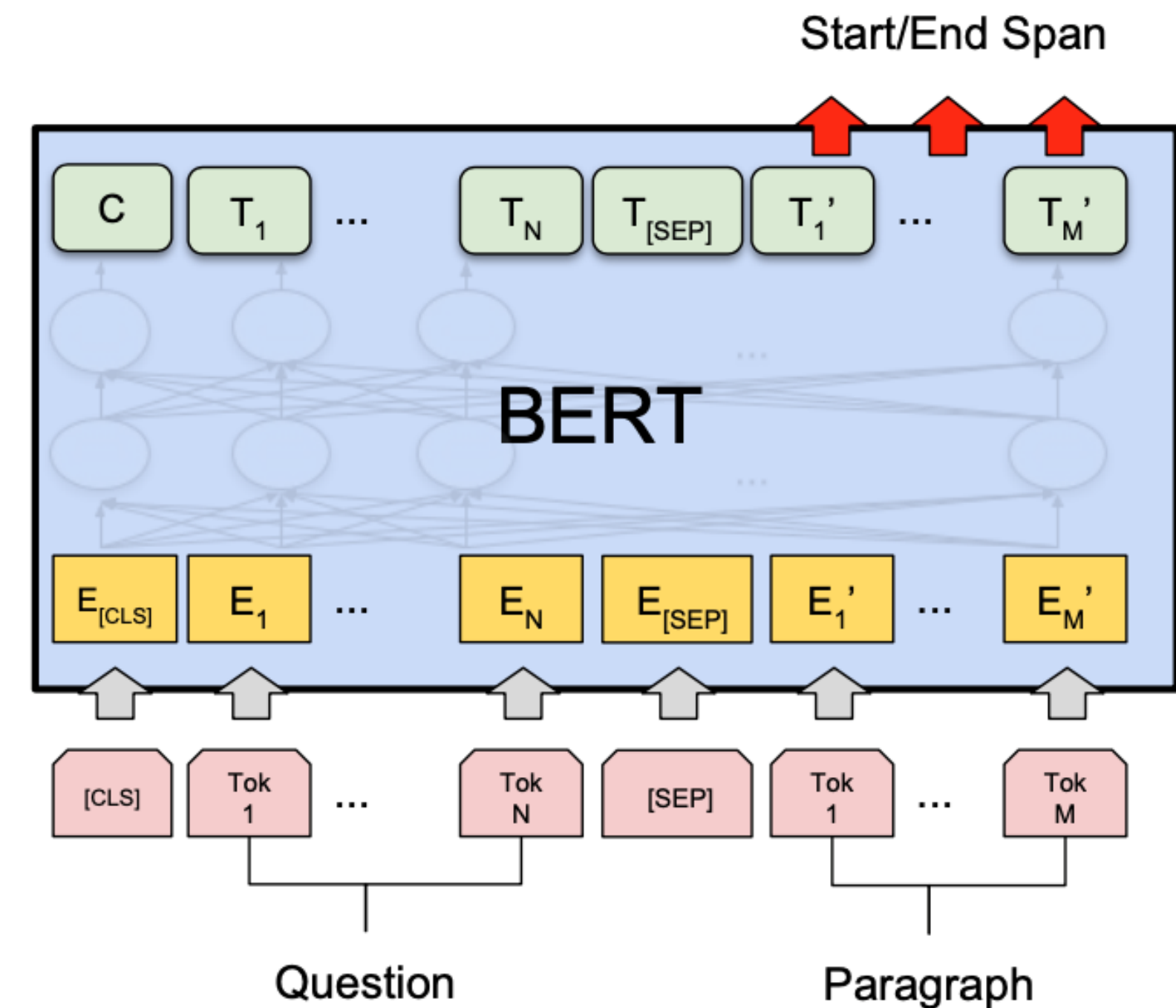where $\mathbf{h}_i$ is the hidden vector of $c_i$, returned by BERT

34

# 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 $\mathbf{h}_{\text{start}}, \mathbf{h}_{\text{end}}$ (e.g., 768 x 2 = 1536) are optimized together for $\mathcal{L}$.

- It works amazing 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)



Start/End Span

BERT

Question          Paragraph

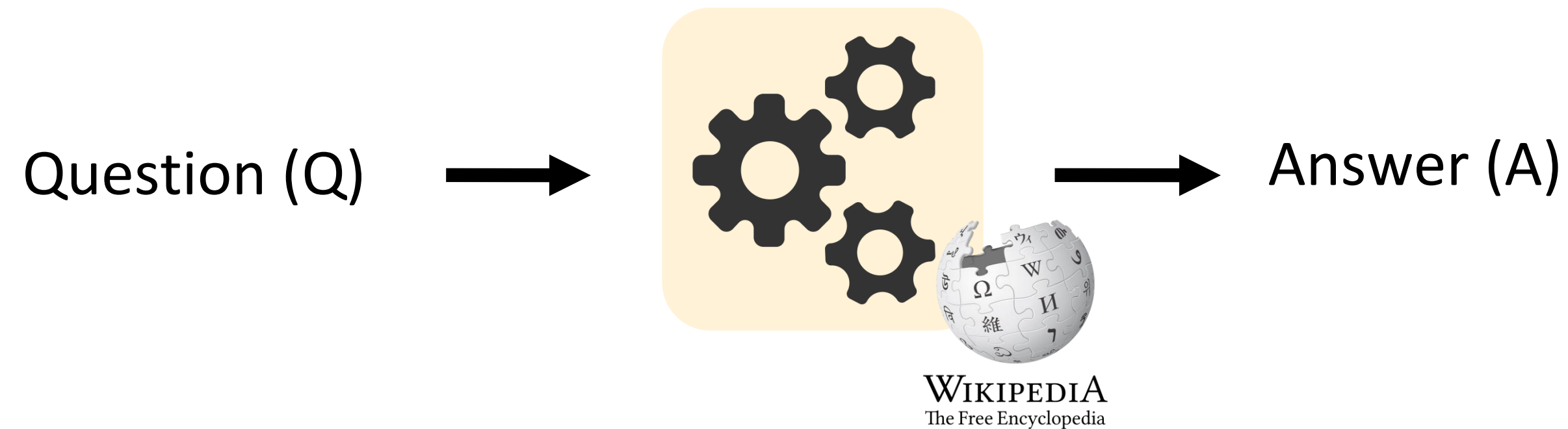# Comparisons between BiDAF and BERT models

- BERT model has many many more parameters (110M or 330M)
  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).
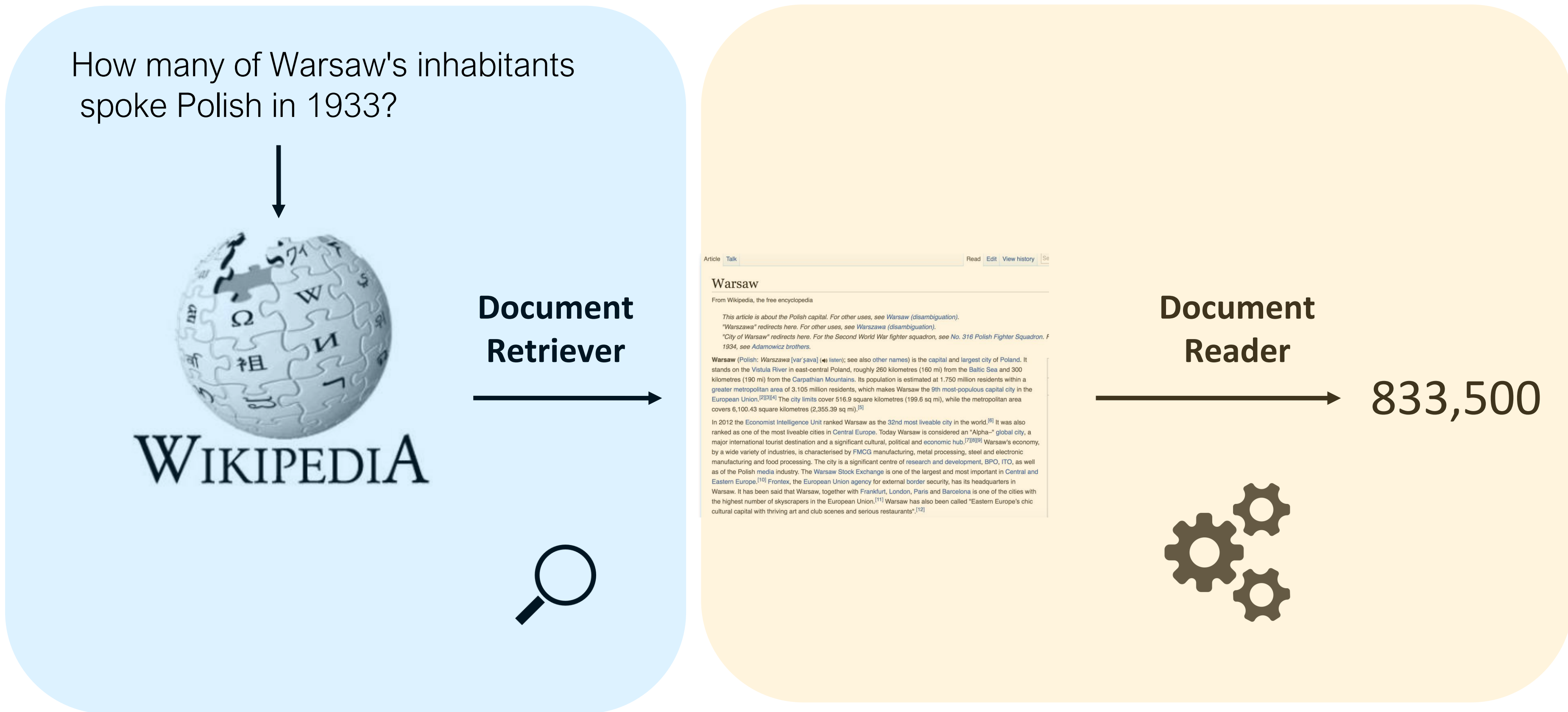
Pre-training is clearly a game changer but it is expensive..

# 3. Open-domain question answering



Question (Q) ➡️ ⚙️ ➡️ Answer (A)

WIKIPEDIA
The Free Encyclopedia

- 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 and a more practical problem!

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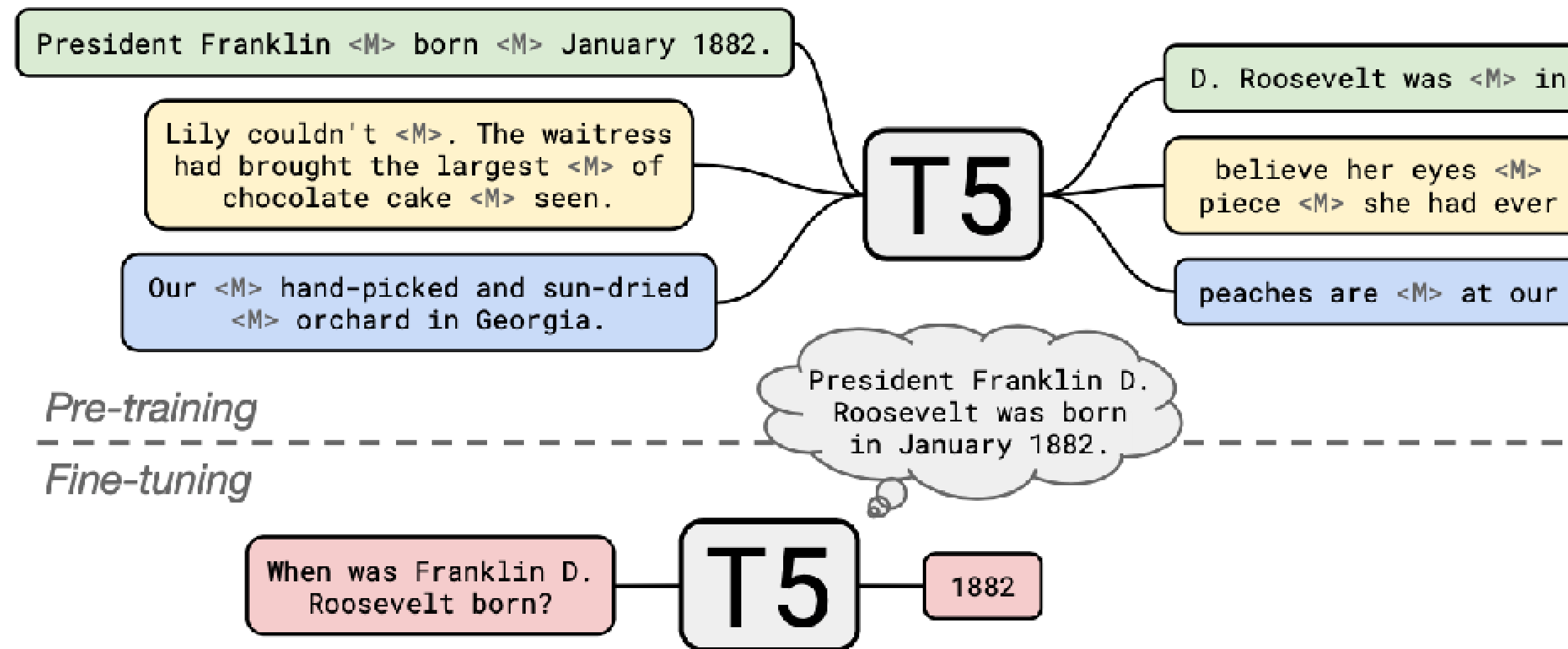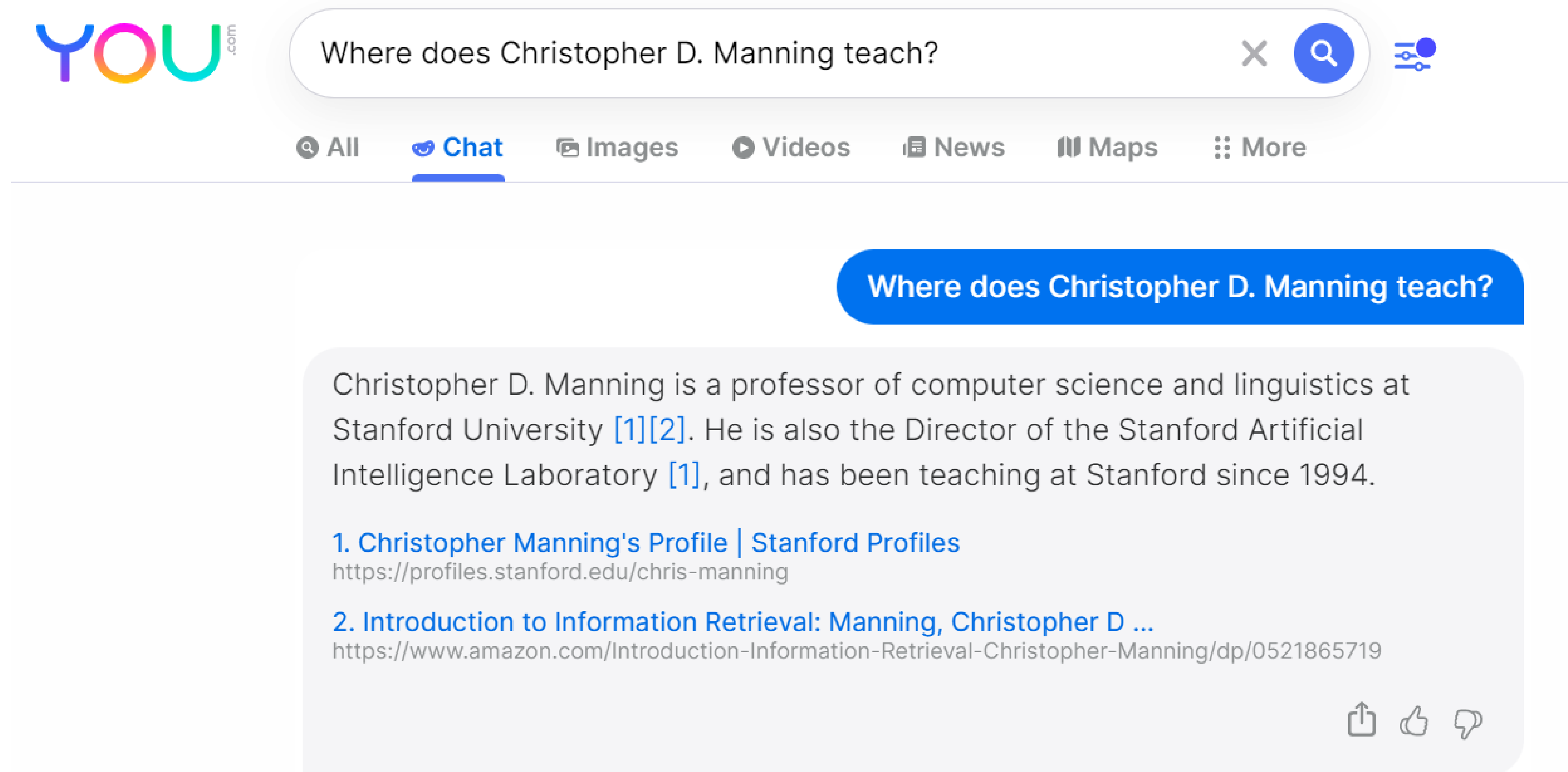
# Retriever-reader framework



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

**Document Retriever**

**Document Reader**

833,500

Chen et al., 2017. Reading Wikipedia to Answer Open-domain Questions

# Large language models can do open-domain QA well

- … without an explicit retriever stage



Roberts et al., 2020. How Much Knowledge Can You Pack Into the Parameters of a Language Model?

# Large language model-based QA (with web search!)

# Problems with large language model-based QA



**Seems totally reasonable!**

**But (1) it's not his most cited paper, and (2) it doesn't have that many citations. Yikes! Also the reference to a web page doesn't help.**

YOU.com

What is the most cited paper by Christopher D. Manning?    ✕   🔍    ⚙

All    🗨 **Chat**    🖼 Images    ▶ Videos    📰 News    📊 Maps    ⠿ More

**What is the most cited paper by Christopher D. Manning?**

The most cited paper by Christopher D. Manning is "Effective Approaches to Attention-Based Neural Machine Translation", which was co-authored by Minh-Thang Luong [1], Hieu Pham, and Christopher D. Manning. This paper has been cited over 18,400 times and is one of the most influential papers in the field of Natural Language Processing.

1. Effective Approaches to Attention-based Neural Machine Translation
https://arxiv.org/abs/1508.04025

🖤 👍 👎

Ask me anything...