

AI VIETNAM
All-in-One Course
(TA Session)

End-to-end Question Answering

Project – P2



AI VIET NAM
[@aivietnam.edu.vn](http://aivietnam.edu.vn)

Dinh-Thang Duong - TA

Outline

- Introduction
- Module: Reader
- Module: Retriever
- Question

Introduction

Introduction

❖ Getting Started



Most famous
search
engines

- 百度热搜 > 换一换
- ↑ 中国网络文明大会将有这些安排
 - 3 抓好抗高温热害干旱夺秋粮丰收
 - 新郎拍婚纱照时遭雷击 人已去世 热
 - 4 怪鱼锁定！专业人员下洞捉拿
 - 芬兰女总理哽咽：我是人 也想找... 热
 - 5 “我现在都在猪圈休息”



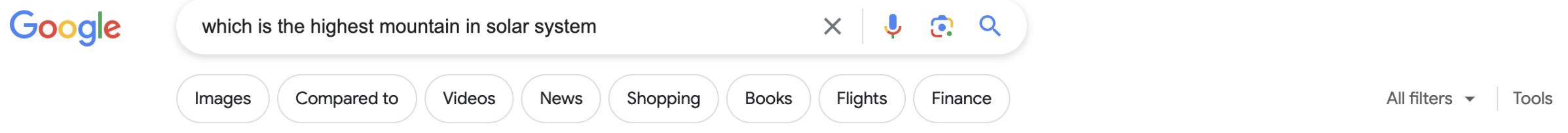
Yandex

Finds everything

Search

Introduction

❖ Getting Started



The image shows a Google search results page. The search bar at the top contains the query "which is the highest mountain in solar system". Below the search bar are several filter buttons: Images, Compared to, Videos, News, Shopping, Books, Flights, and Finance. On the right side of the search bar are icons for microphone, camera, and search. At the bottom right of the search bar are links for "All filters" and "Tools".

Olympus Mons

The highest mountain and volcano in the Solar System is on the planet Mars. It is called **Olympus Mons** and is 16 miles (24 kilometers) high which makes it about three times higher than Mt. Everest.



Cool Cosmos

<https://coolcosmos.ipac.caltech.edu> › ask › 199-Where-is... · · ·

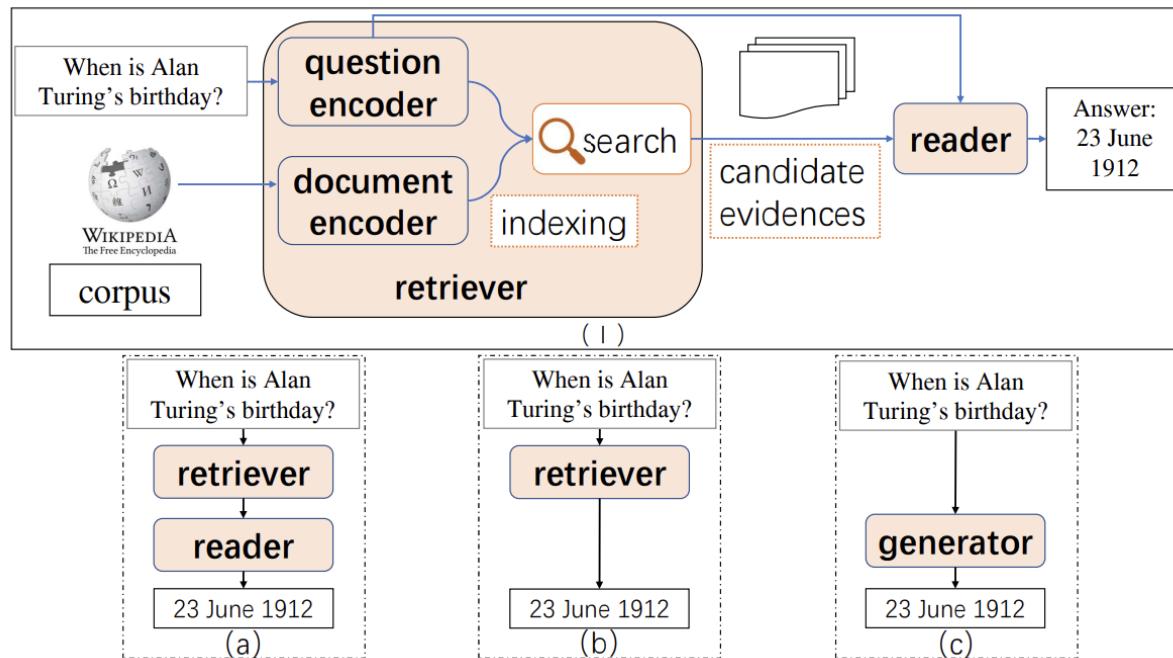
Where is the highest mountain in our Solar System?



Answer from Google for the query question: Open-domain Question Answering

Introduction

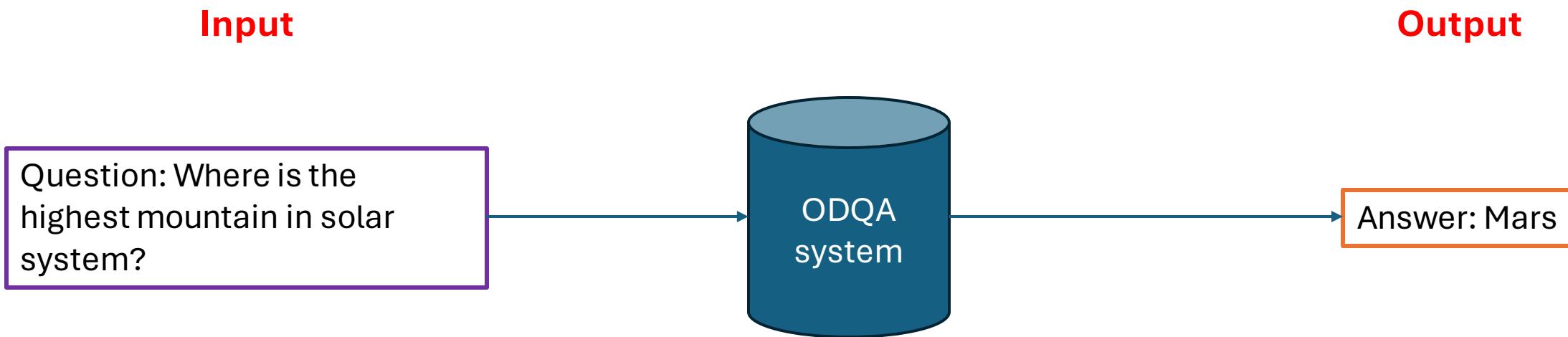
❖ What is Open-domain Question Answering?



Open-domain Question Answering (ODQA): Refers to the task of providing answers to questions over a broad range of topics, leveraging vast corpora of unstructured text without being restricted to a specific domain.

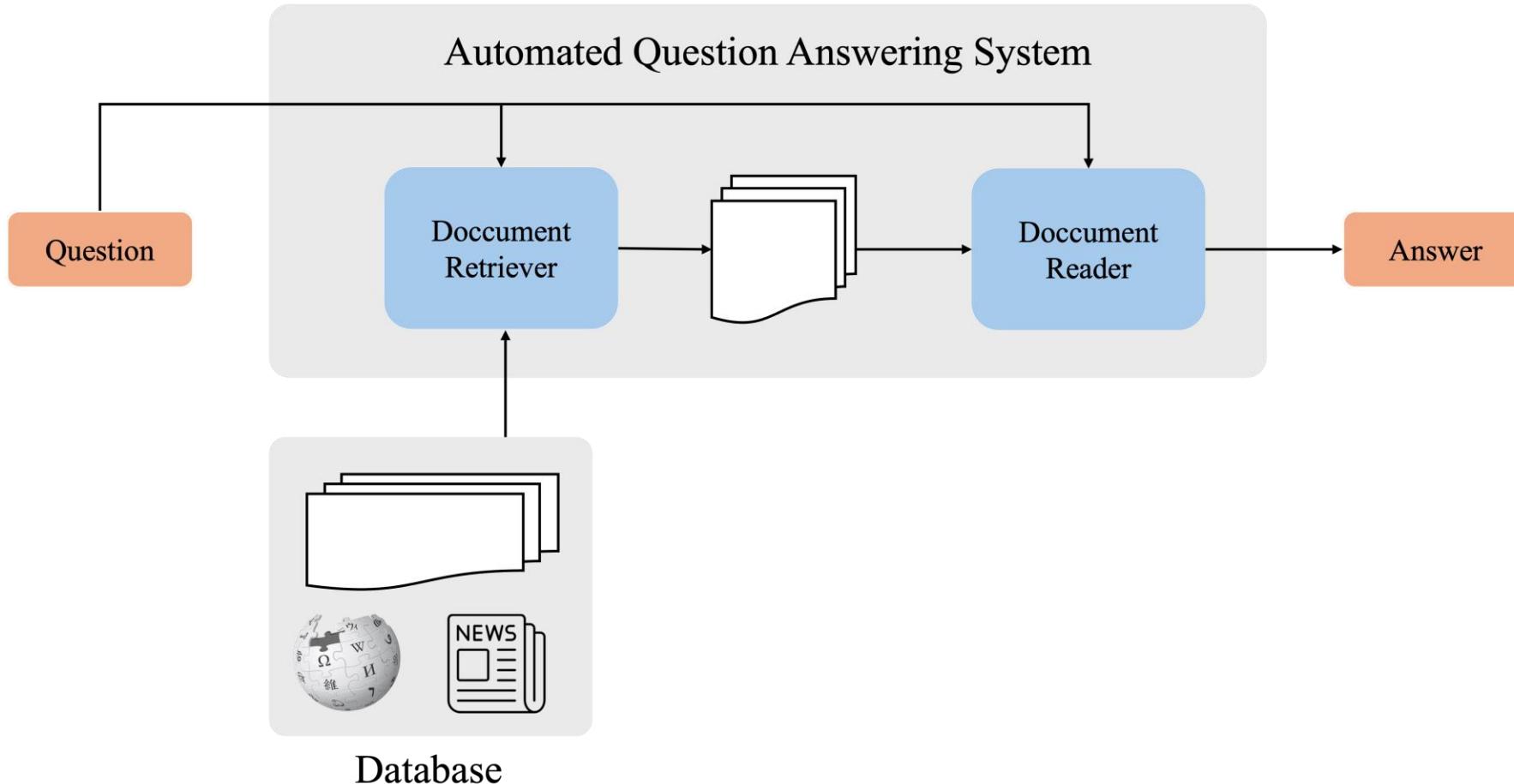
Introduction

❖ Open-domain Question Answering I/O



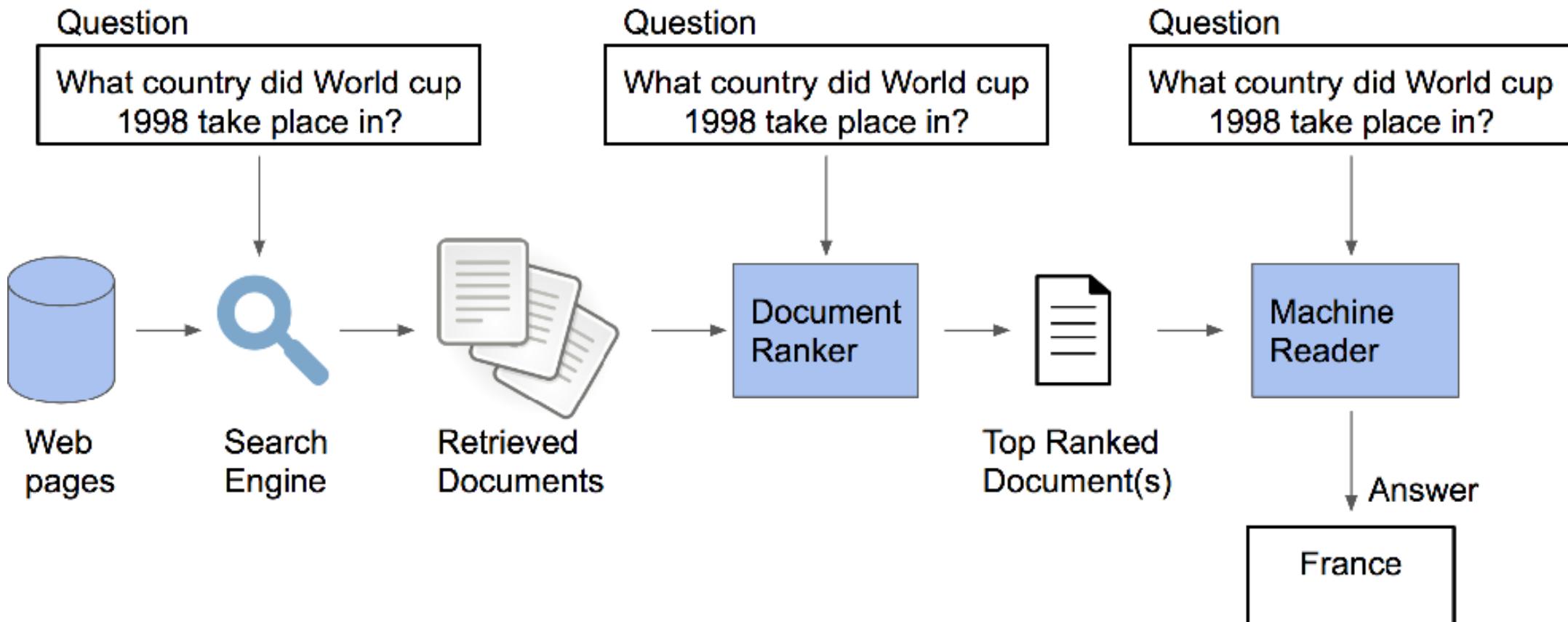
Introduction

❖ End-to-end Question Answering System I/O



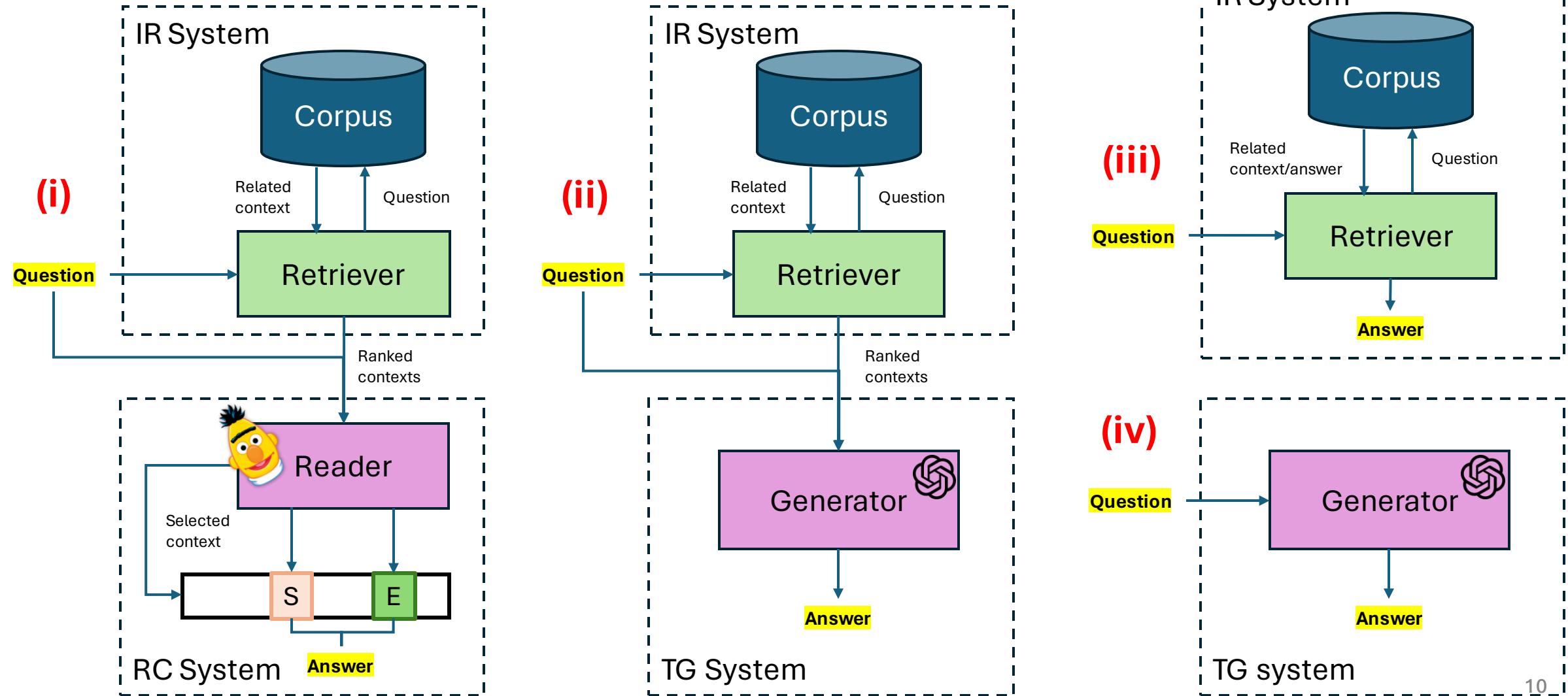
Introduction

❖ Open-domain Question Answering I/O



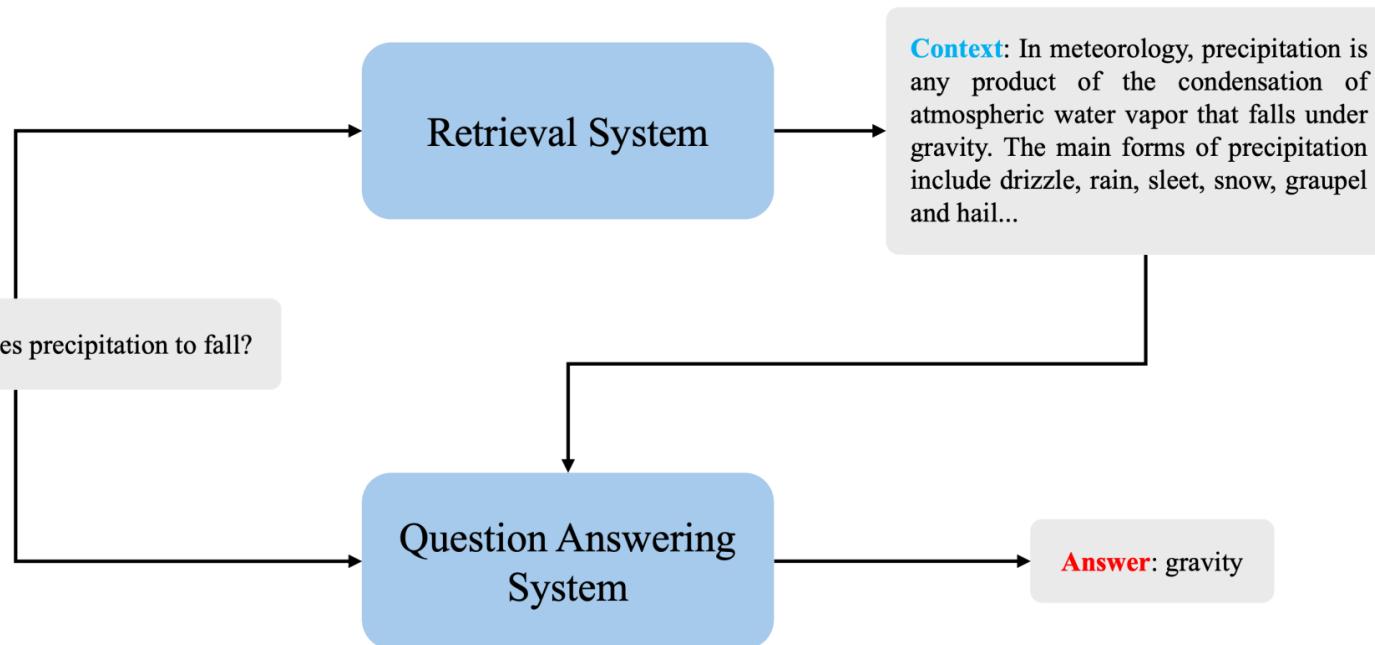
Introduction

❖ Open-domain Question Answering Categories



Introduction

❖ End-to-end Question Answering System



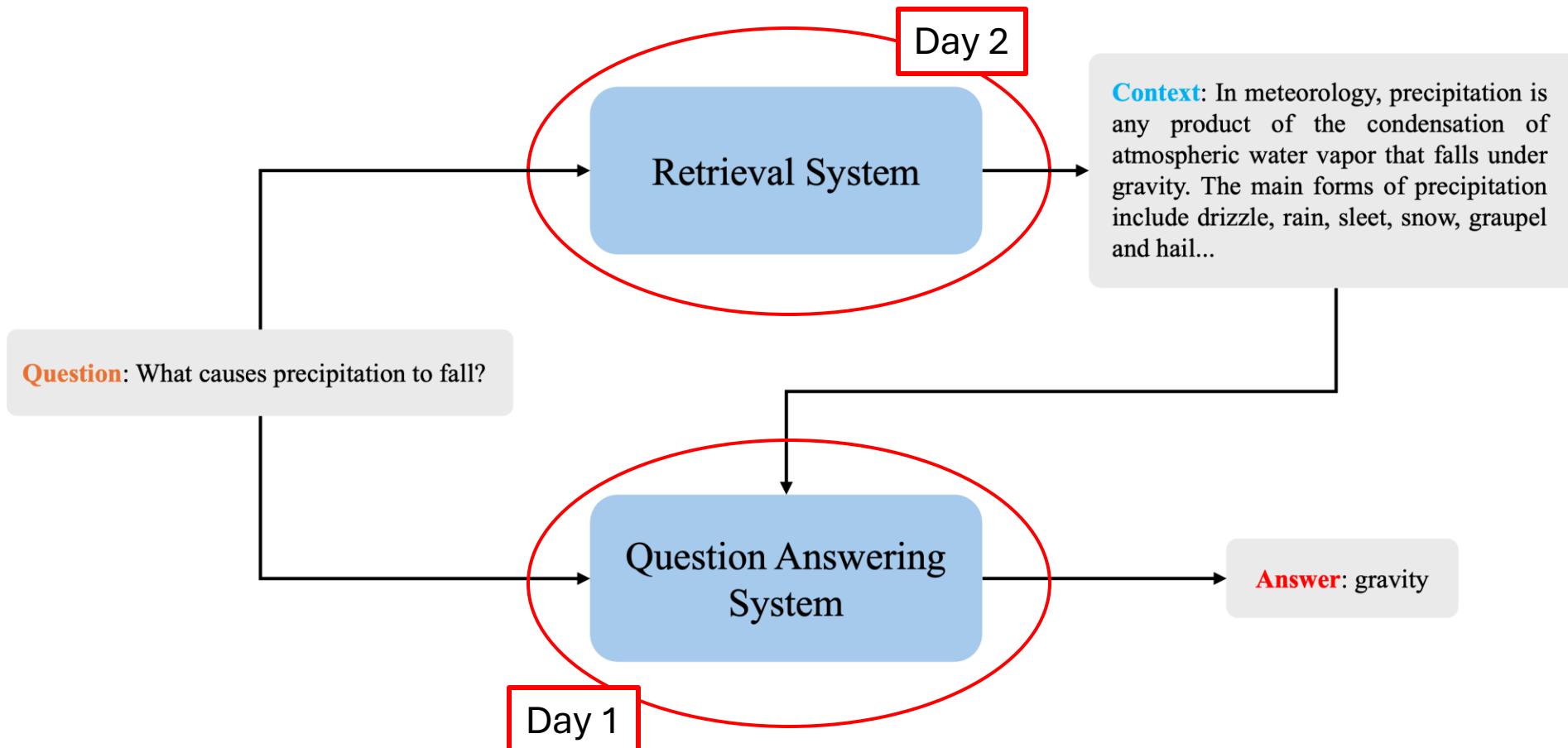
End-to-end QA system consists of:

1. Retriever
2. Reader

Introduction

❖ Project Description

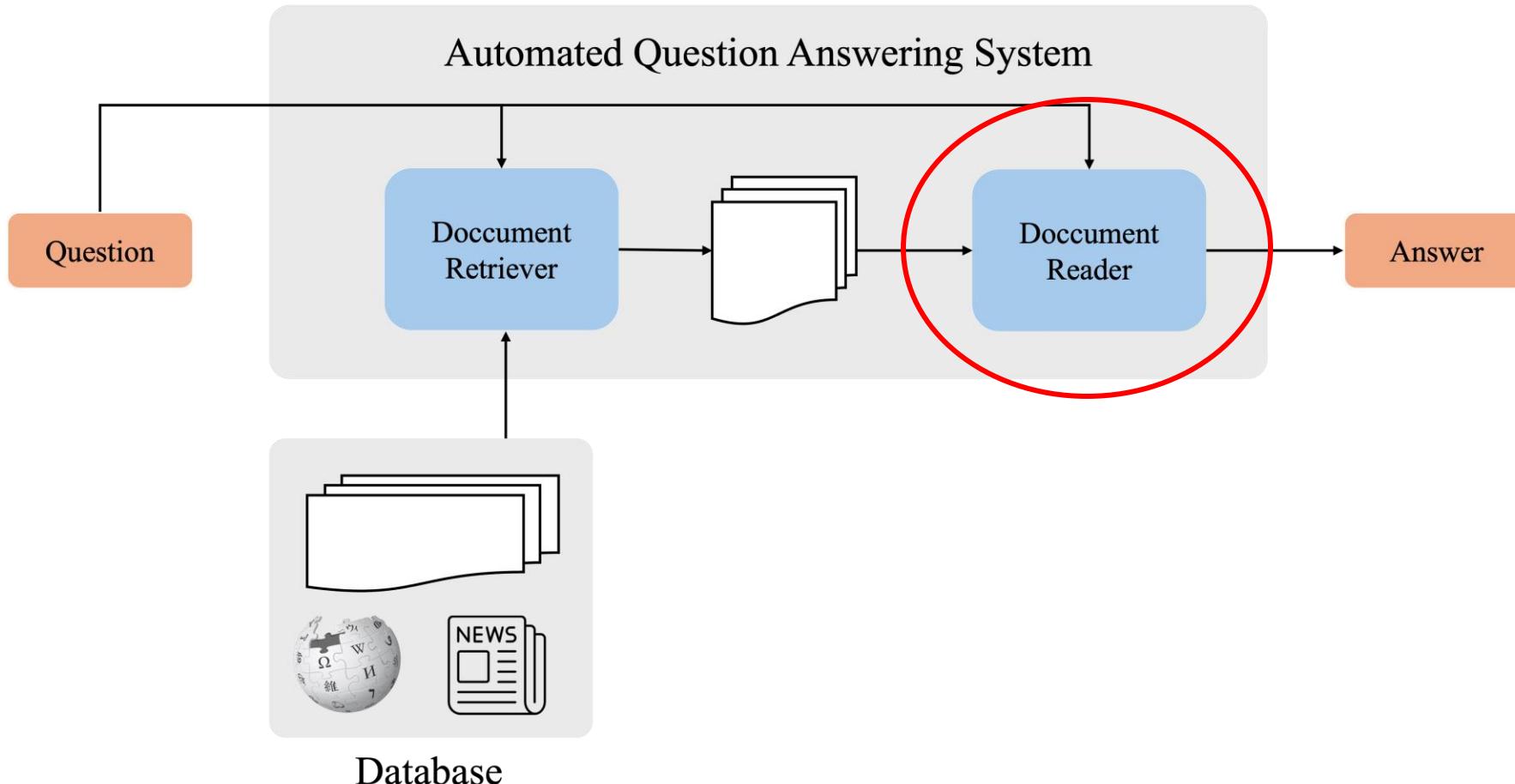
Description: In this project, we will build an end-to-end QA system that can retrieve relevant contexts, extract best answer for a given question. The domain of the question can be varied in different topics.



Module: Reader

Module: Reader

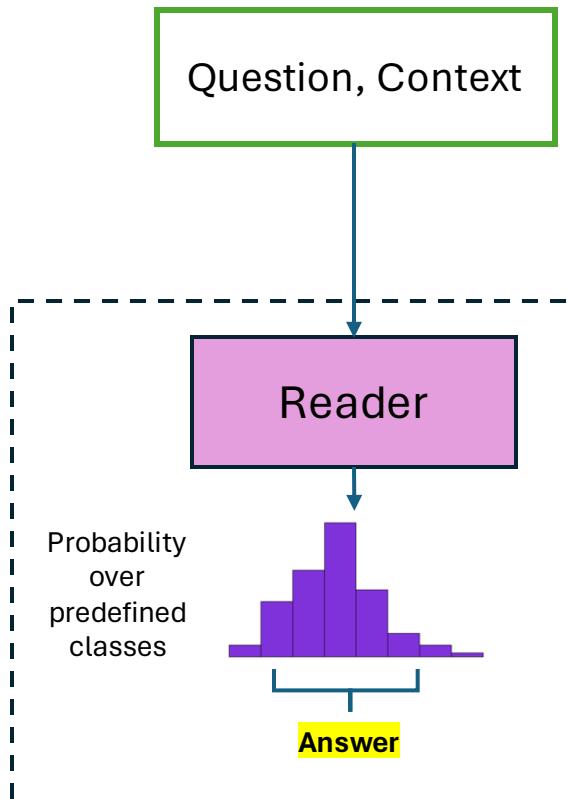
❖ E2E QA Pipeline



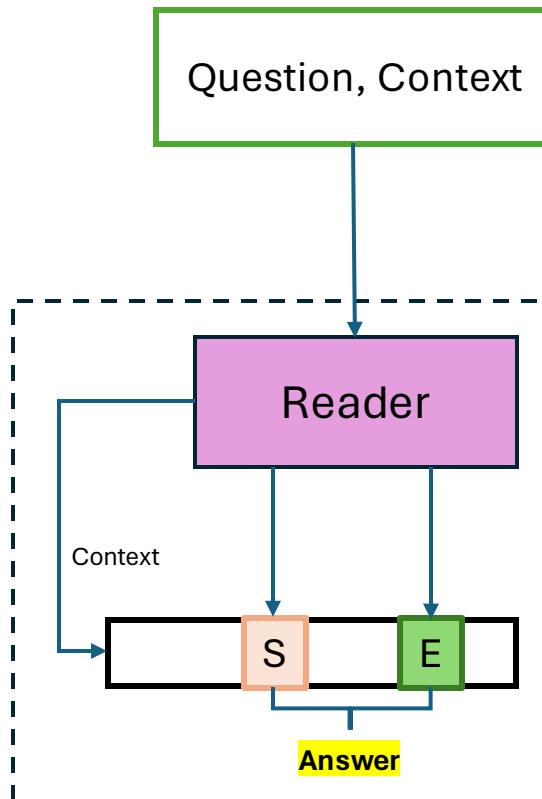
Module: Reader

❖ Type of Reader

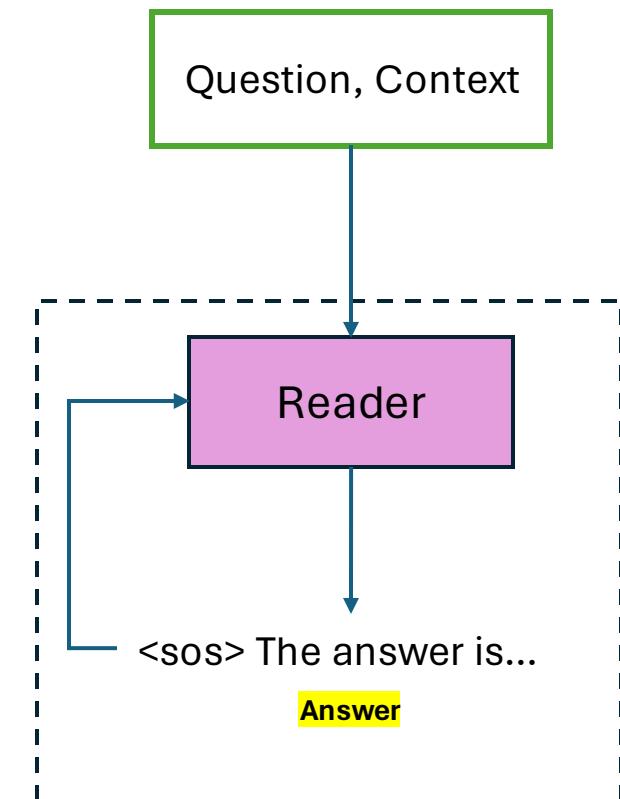
I. Classification Approach



II. Extractive Approach



III. Generative Approach



Module: Reader

❖ Reader: Extractive Extraction

Input

Question: Where is the highest mountain in solar system?

Output

QA
(Reader)

Context: The highest mountain and volcano in the Solar System is on the planet Mars. It is called Olympus Mons and is 16 miles (24 kilometers) high which makes it about three times higher than Mt. Everest.

- Start positions
- End positions

Module: Reader

❖ Model Design

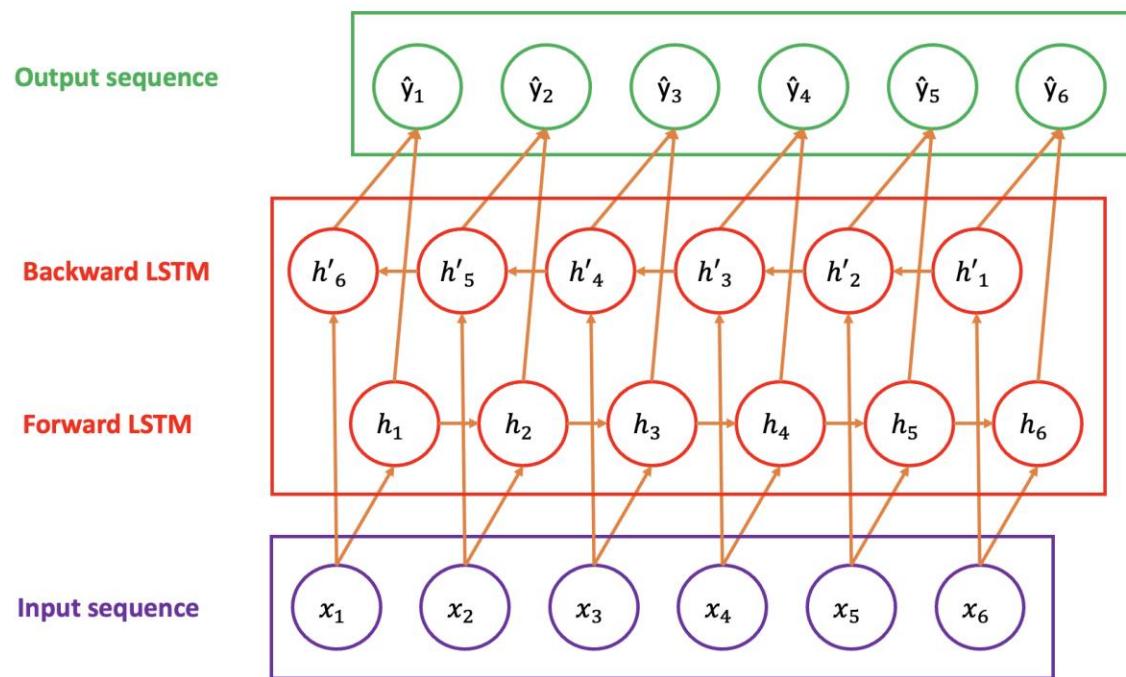
```
1 INPUT_QUESTION = 'What is my name?'
2 INPUT_CONTEXT = 'My name is AI Vietnam and I live in Vietnam.'
3 pipe(question=INPUT_QUESTION, context=INPUT_CONTEXT)
```

```
{'score': 0.97179114818573, 'start': 11, 'end': 21, 'answer': 'AI Vietnam'}
```

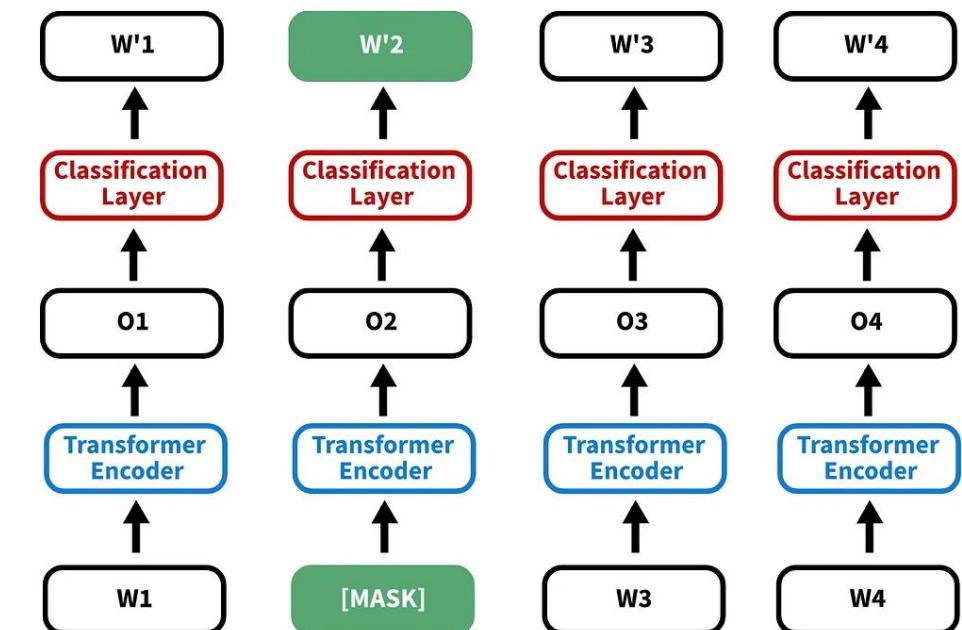
The answer span starts from index 11 to 21 in context text

Module: Reader

❖ Model Design: LSTM and Transformer Encoder



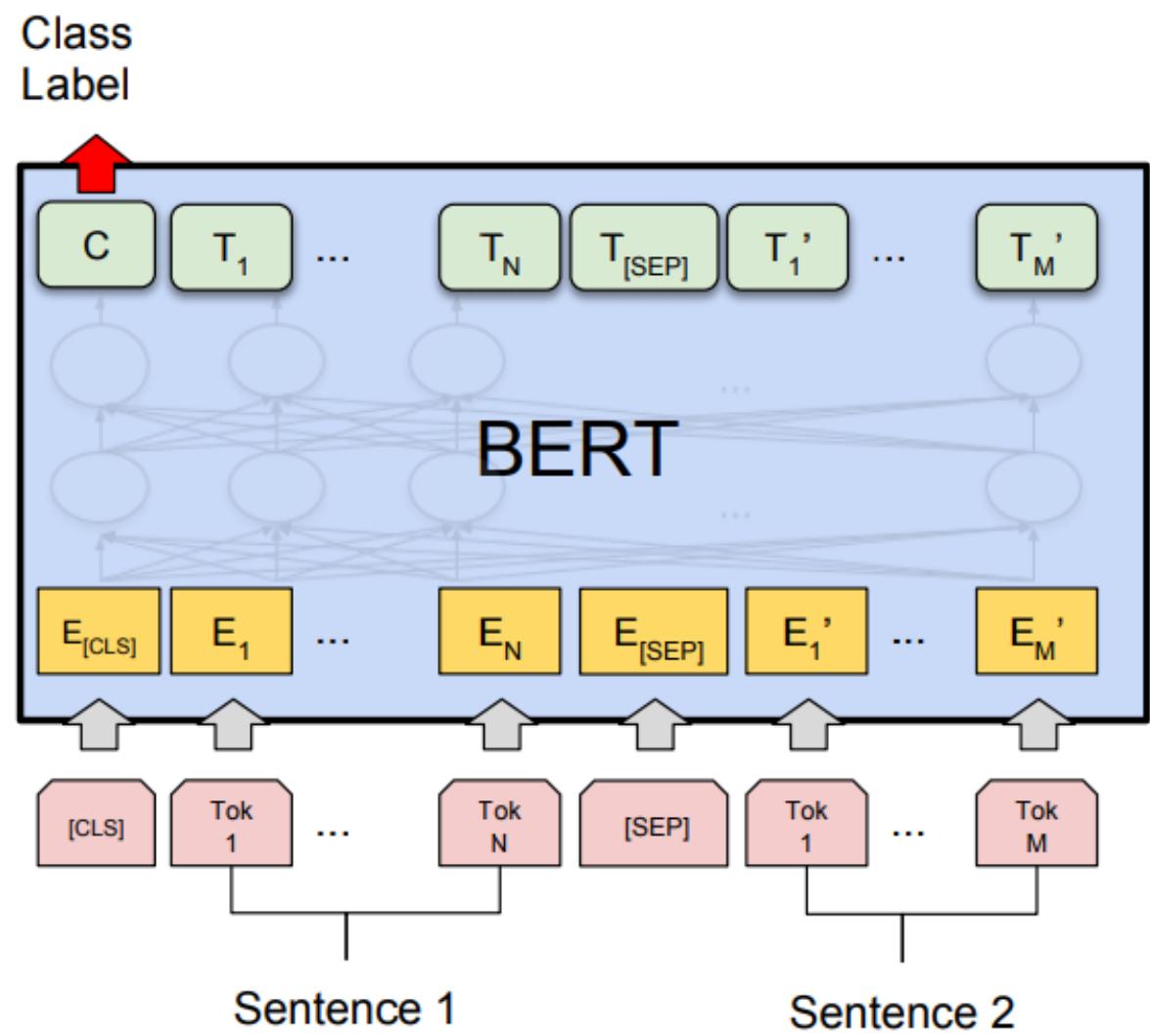
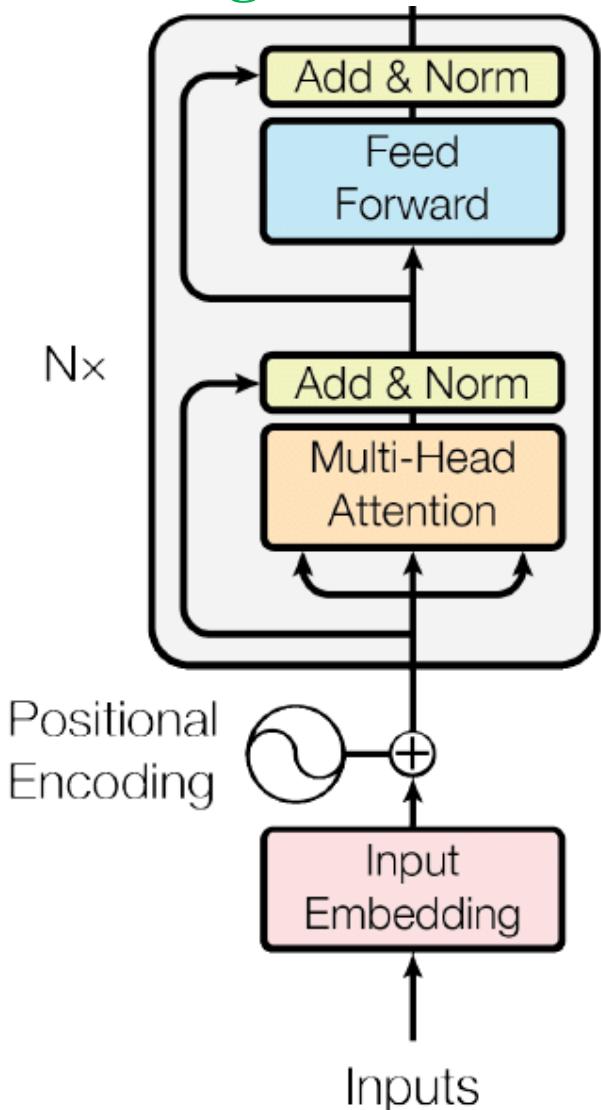
LSTM-based



Transformer Encoder-based

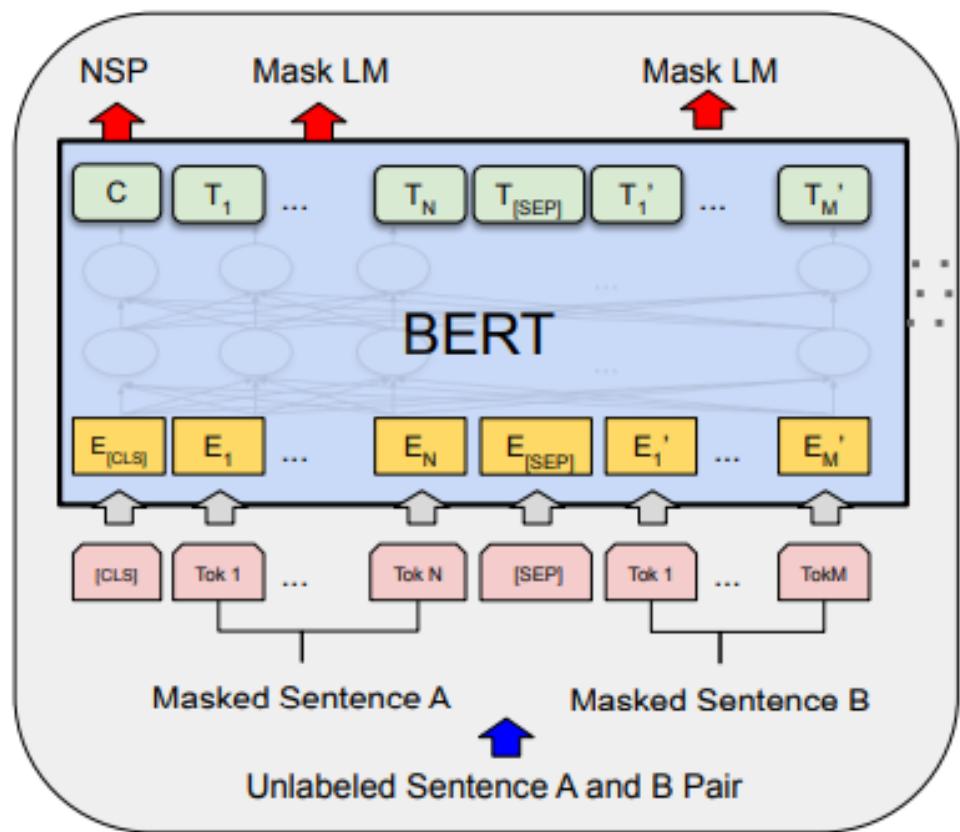
Module: Reader

❖ Module Design: BERT

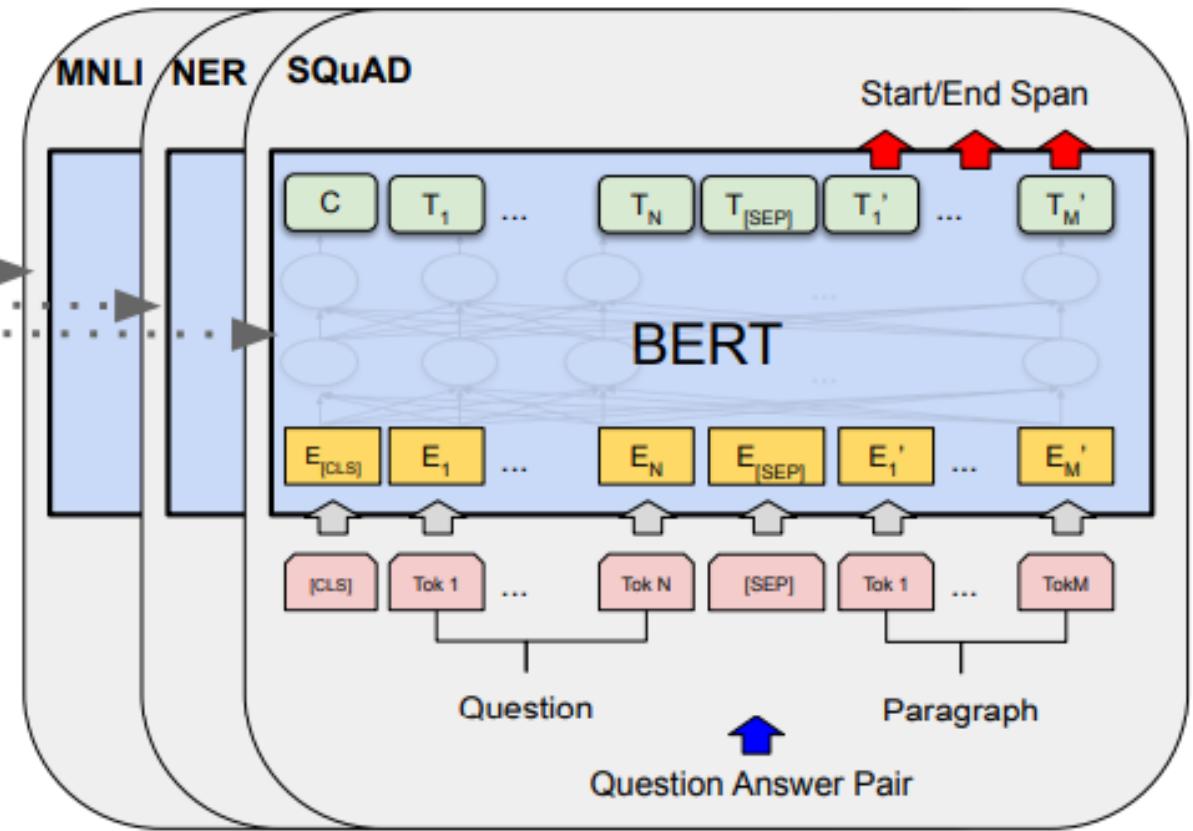


Module: Reader

❖ From original BERT to BERT for QA



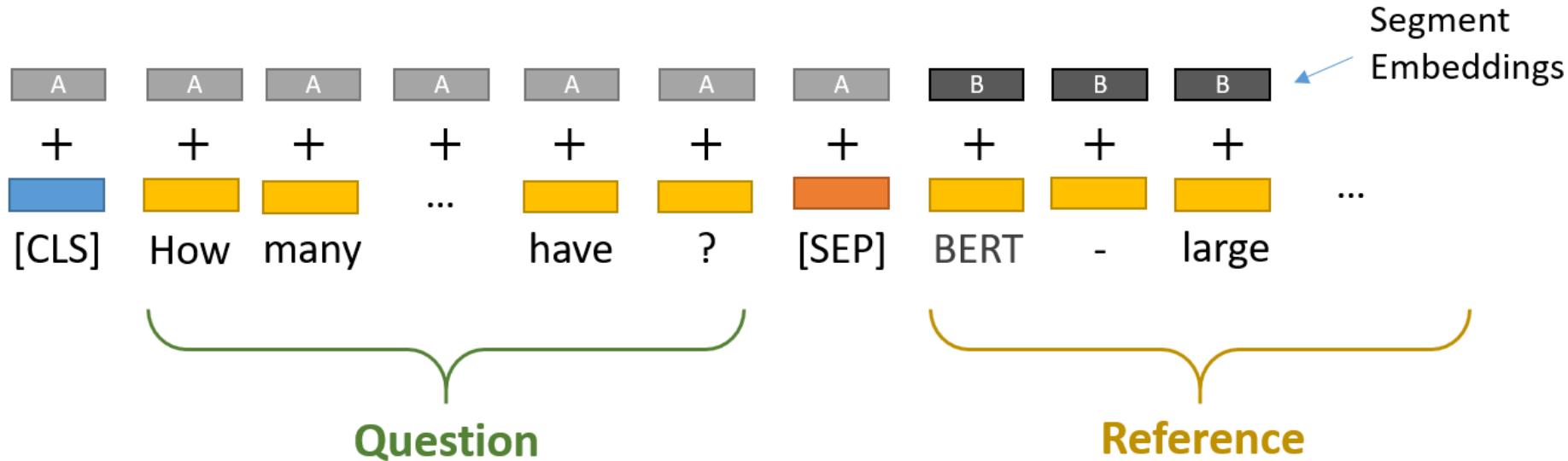
Pre-training



Fine-Tuning

Module: Reader

❖ BERT Input for QA

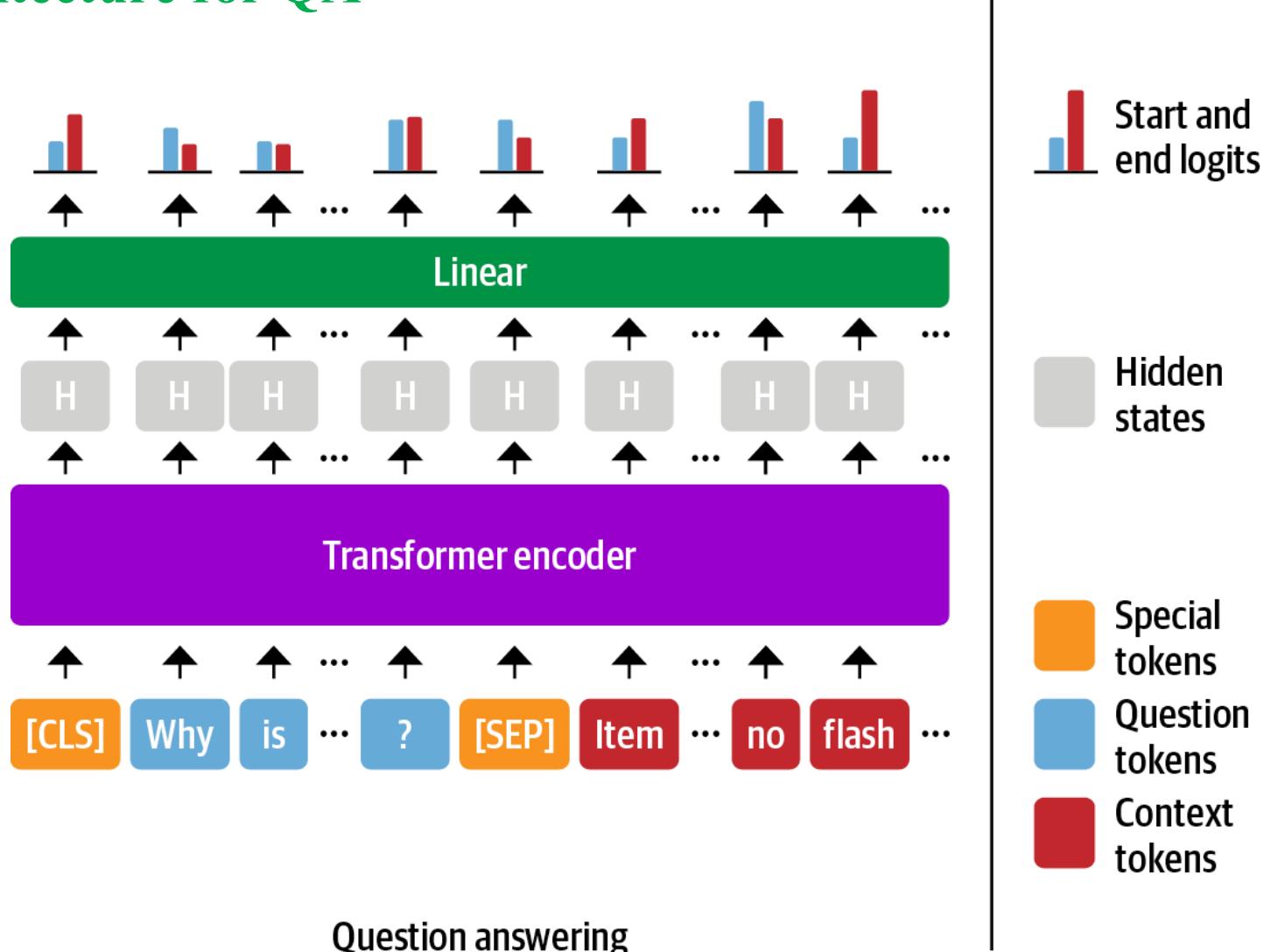


Question: How many parameters does BERT-large have?

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.

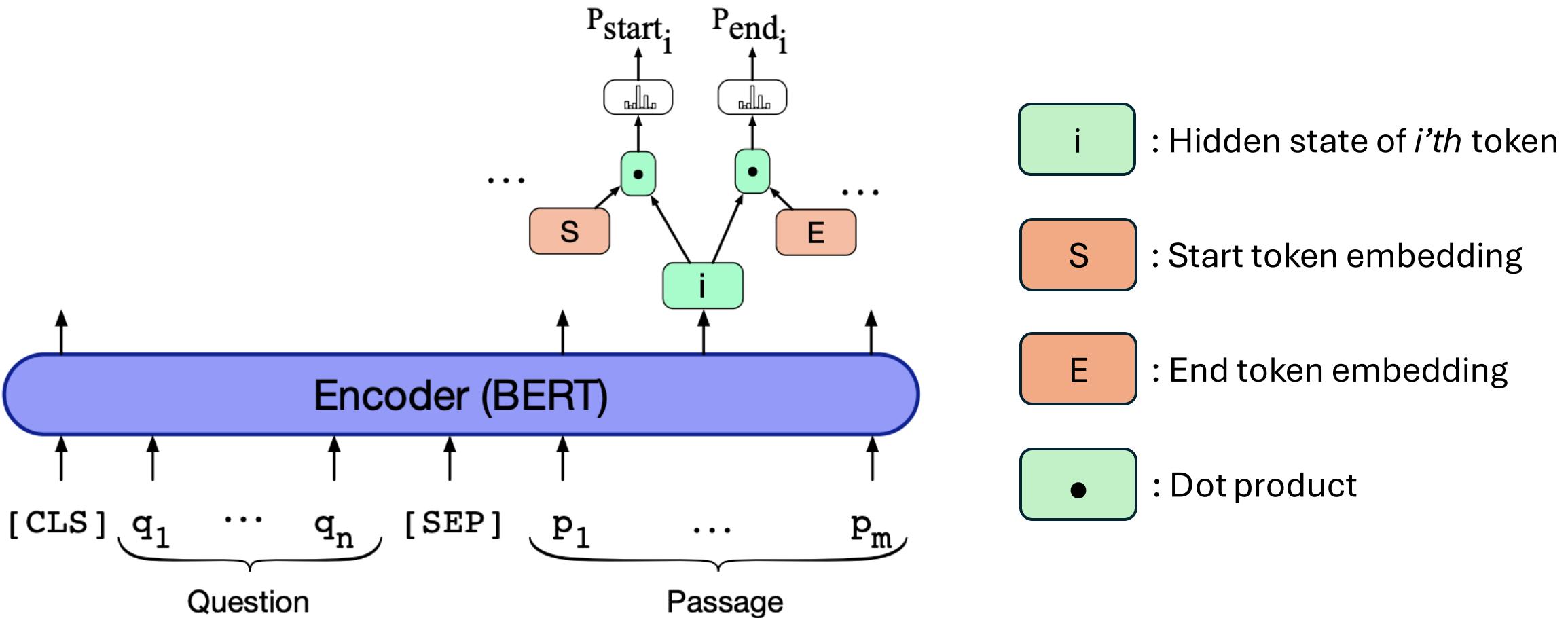
Module: Reader

❖ BERT Architecture for QA



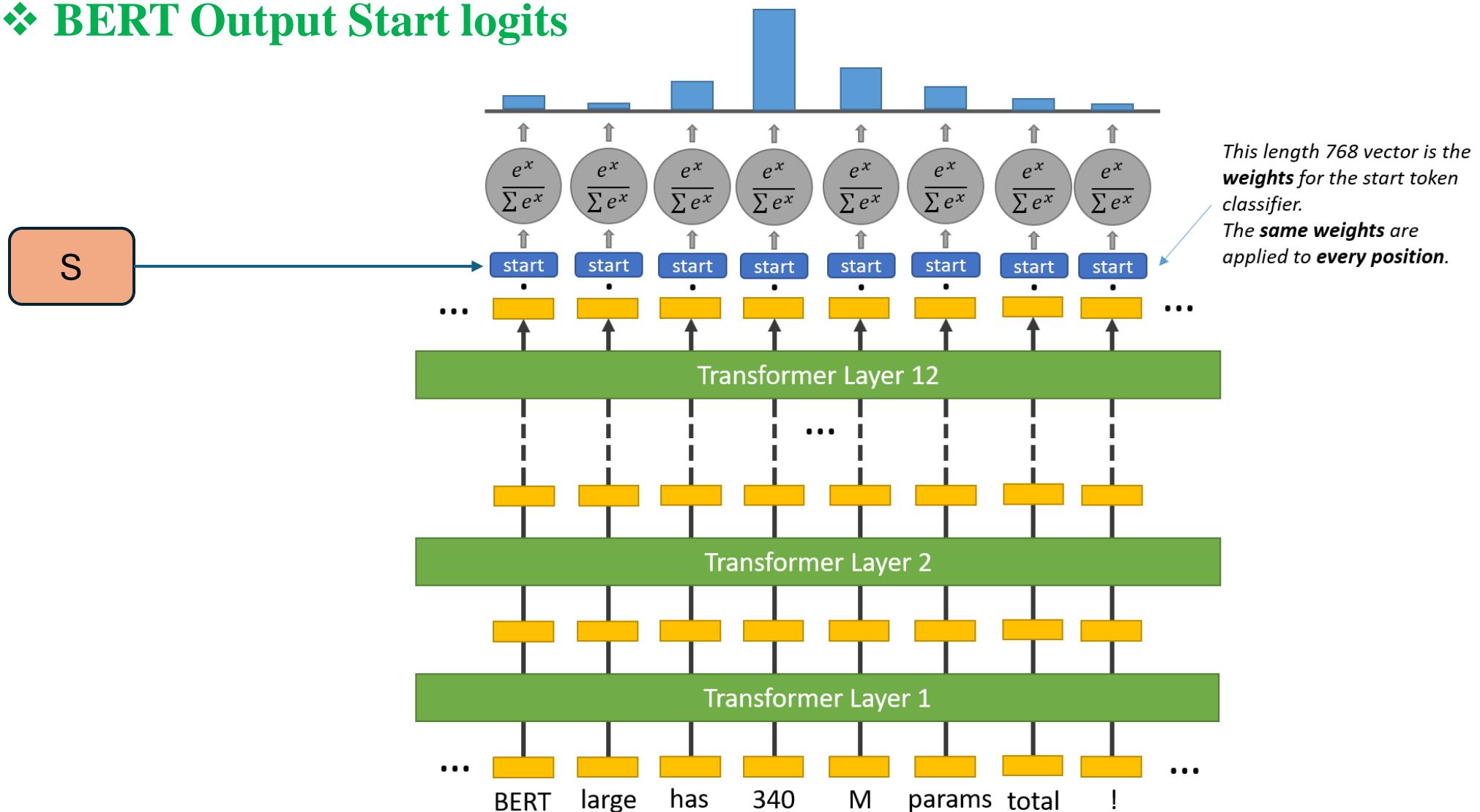
Module: Reader

❖ BERT Architecture for QA



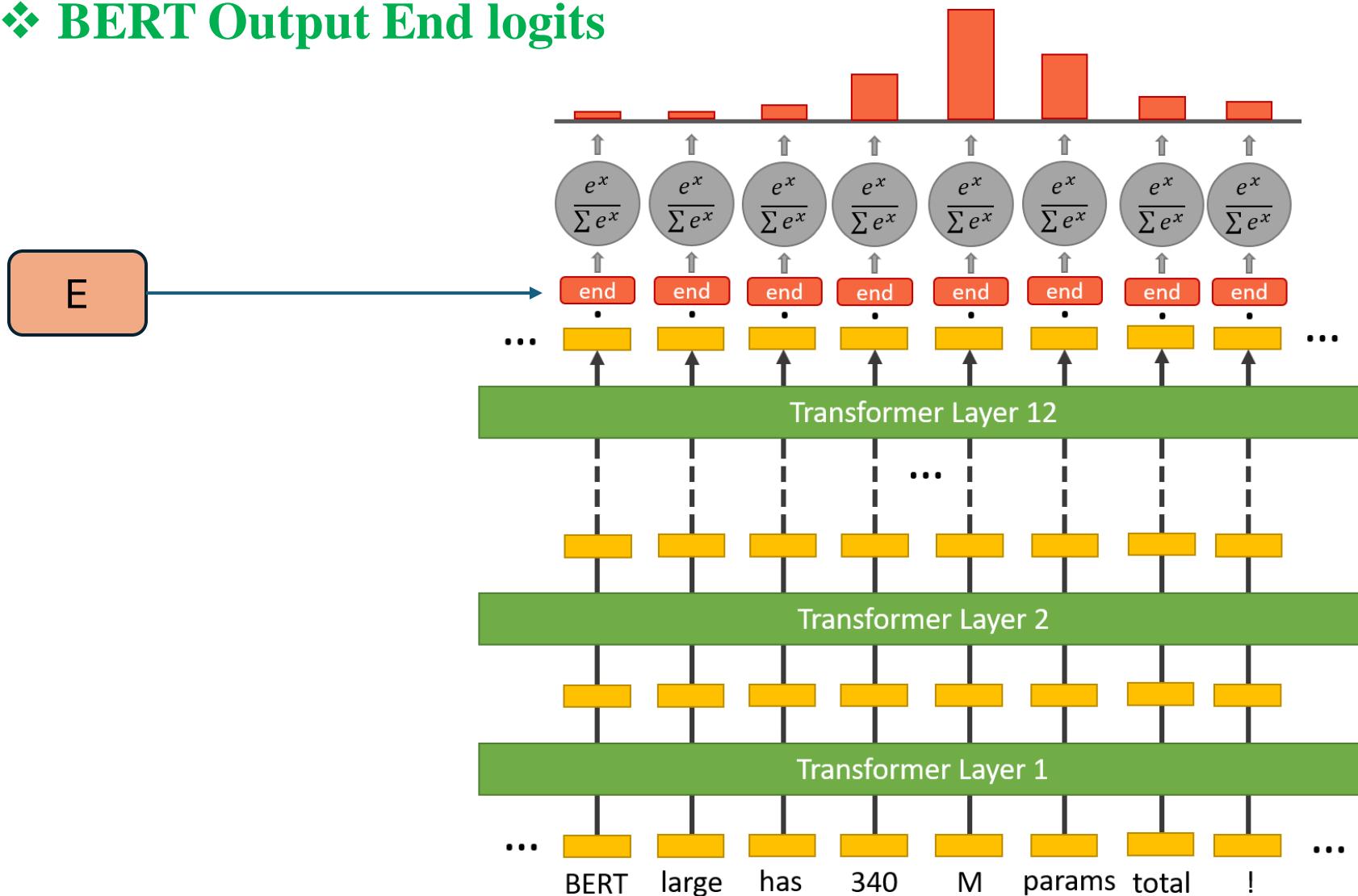
Module: Reader

❖ BERT Output Start logits



Module: Reader

❖ BERT Output End logits



Module: Reader

❖ Step 1: Install and import libraries

```
1 !pip install -qq datasets==2.16.1 evaluate==0.4.1 transformers[sentencepiece]==4.35.2
2 !pip install -qq accelerate==0.26.1
3 !apt install git-lfs
```

```
----- 507.1/507.1 kB 4.3 MB/s eta 0:00:00
----- 84.1/84.1 kB 13.4 MB/s eta 0:00:00
----- 115.3/115.3 kB 18.8 MB/s eta 0:00:00
----- 134.8/134.8 kB 20.9 MB/s eta 0:00:00
----- 134.8/134.8 kB 22.2 MB/s eta 0:00:00
----- 270.9/270.9 kB 7.1 MB/s eta 0:00:00
```

Reading package lists... Done

Building dependency tree... Done

Reading state information... Done

git-lfs is already the newest version (3.0.2-1ubuntu0.2).

0 upgraded, 0 newly installed, 0 to remove and 33 not upgraded.

Module: Reader

❖ Step 1: Install and import libraries

```
1 import numpy as np
2 from tqdm.auto import tqdm
3 import collections
4
5 import torch
6
7 from datasets import load_dataset
8 from transformers import AutoTokenizer
9 from transformers import AutoModelForQuestionAnswering
10 from transformers import TrainingArguments
11 from transformers import Trainer
12 import evaluate
13
14 device = torch.device("cuda") if torch.cuda.is_available() else torch.device("cpu")
```



Module: Reader

❖ Step 2: Setup config

```
1 # Sử dụng mô hình "distilbert-base-uncased"
2 # làm mô hình checkpoint
3 MODEL_NAME = "distilbert-base-uncased"
4
5 # Độ dài tối đa cho mỗi đoạn văn bản
6 # sau khi được xử lý
7 MAX_LENGTH = 384
8
9 # Khoảng cách giữa các điểm bắt đầu
10 # của các đoạn văn bản liên tiếp
11 STRIDE = 128
```

Module: Reader

❖ Step 3: Download dataset

The screenshot shows the official SQuAD2.0 website. The top banner features the title "SQuAD2.0" and the subtitle "The Stanford Question Answering Dataset".

What is SQuAD?

Stanford Question Answering Dataset (SQuAD) is a reading comprehension dataset, consisting of questions posed by crowdworkers on a set of Wikipedia articles, where the answer to every question is a segment of text, or *span*, from the corresponding reading passage, or the question might be unanswerable.

SQuAD2.0 combines the 100,000 questions in SQuAD1.1 with over 50,000 unanswerable questions written adversarially by crowdworkers to look similar to answerable ones. To do well on SQuAD2.0, systems must not only answer questions when possible, but also determine when no answer is supported by the paragraph and abstain from answering.

[Explore SQuAD2.0 and model predictions](#)

[SQuAD2.0 paper \(Rajpurkar & Jia et al. '18\)](#)

Leaderboard

SQuAD2.0 tests the ability of a system to not only answer reading comprehension questions, but also abstain when presented with a question that cannot be answered based on the provided paragraph.

Rank	Model	EM	F1
	Human Performance <i>Stanford University</i> (Rajpurkar & Jia et al. '18)	86.831	89.452
1 Jun 04, 2021	IE-Net (ensemble) <i>RICOH_SRCB_DML</i>	90.939	93.214
2 Feb 21, 2021	FPNet (ensemble) <i>Ant Service Intelligence Team</i>	90.871	93.183
3 May 16, 2021	IE-NetV2 (ensemble) <i>RICOH_SRCB_DML</i>	90.860	93.100
4 Apr 06, 2020	SA-Net on Albert (ensemble) <i>QIANXIN</i>	90.724	93.011

SQuAD v2 dataset

Module: Reader

❖ Step 3: Download dataset

SQuAD 2.0	
Train	
Total examples	130,319
Negative examples	43,498
Total articles	442
Articles with negatives	285
Development	
Total examples	11,873
Negative examples	5,945
Total articles	35
Articles with negatives	35
Test	
Total examples	8,862
Negative examples	4,332
Total articles	28
Articles with negatives	28

Answer type	Percentage	Example
Date	8.9%	19 October 2023
Other Numeric	10.9%	12
Person	12.9%	Thomas Coke
Location	4.4%	Germany
Other Entity	15.3%	ABC Sports
Common Noun Phrase	31.8%	property damage
Adjective Phrase	3.9%	second-largest
Verb Phrase	5.5%	returned to Earth
Clause	3.7%	to avoid trivialization
Other	2.7%	quietly

Module: Reader

❖ Step 3: Download dataset

Context: The English name "Normans" comes from the French words Normans/Normanç, plural of Normant, modern French normand, which is itself borrowed from (Old Low Franconian Nortmann "Northman" or directly from Old Norse Nordmaðr, Latinized variously as Nortmannus, Normannus, or Nordmannus (recorded in Medieval Latin, **9th century**) to mean "Norseman, **Viking**".

Q1: What causes precipitation to fall?

A1: Viking

Q2: When was the Latin version of the word Norman first recorded?

A2: 9th century

Q3: When was the French version of the word Norman first recorded?

A3: No Answers

Context: 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”.

Q1: What causes precipitation to fall?

A1: Gravity

Q2: What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?

A2: Graupel

Q3: Where do water droplets collide with ice crystals to form precipitation?

A3: Within a cloud

SQuAD v2 dataset samples

Module: Reader

❖ Step 3: Download dataset

id string · lengths	title string · lengths	context string · lengths	question string · lengths	answers sequence
24	24	3 151	1	
56be85543aeaaa14008c9063	Beyoncé	Beyoncé Giselle Knowles-Carter (/bi...	When did Beyonce start becoming...	{ "text": ["in the late 1990s"]...}
56be85543aeaaa14008c9065	Beyoncé	Beyoncé Giselle Knowles-Carter (/bi...	What areas did Beyonce compete in...	{ "text": ["singing and..."]...}
56be85543aeaaa14008c9066	Beyoncé	Beyoncé Giselle Knowles-Carter (/bi...	When did Beyonce leave Destiny's...	{ "text": ["2003"], "answer_start"... }
56bf6b0f3aeaaa14008c9601	Beyoncé	Beyoncé Giselle Knowles-Carter (/bi...	In what city and state did Beyonce...	{ "text": ["Houston, Texas"]... }
56bf6b0f3aeaaa14008c9602	Beyoncé	Beyoncé Giselle Knowles-Carter (/bi...	In which decade did Beyonce become...	{ "text": ["late 1990s"]... }
56bf6b0f3aeaaa14008c9603	Beyoncé	Beyoncé Giselle Knowles-Carter (/bi...	In what R&B group was she the lead...	{ "text": ["Destiny's Child"]... }

Download SQuAD v2 dataset in [here](#)

Module: Reader

❖ Step 3: Download dataset

```
1 # Download squad dataset từ HuggingFace
2 DATASET_NAME = 'squad_v2'
3 raw_datasets = load_dataset(DATASET_NAME)
```

```
/usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_token.py:88: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens)
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
warnings.warn(
Downloading readme: 100%  8.18k/8.18k [00:00<00:00, 551kB/s]
Downloading data: 100%  16.4M/16.4M [00:00<00:00, 19.5MB/s]
Downloading data: 100%  1.35M/1.35M [00:00<00:00, 4.66MB/s]
Generating train split: 100%  130319/130319 [00:00<00:00, 353109.21 examples/s]
Generating validation split: 100%  11873/11873 [00:00<00:00, 225130.97 examples/s]
```

Module: Reader

❖ Step 3: Download dataset

```
1 raw_datasets
```

```
DatasetDict({  
    train: Dataset({  
        features: ['id', 'title', 'context', 'question', 'answers'],  
        num_rows: 130319  
    })  
    validation: Dataset({  
        features: ['id', 'title', 'context', 'question', 'answers'],  
        num_rows: 11873  
    })  
})
```

```
1 # Print các thông tin Context, Question, và Answer trong dataset  
2 print("Context: ", raw_datasets["train"][0]["context"])  
3 print("Question: ", raw_datasets["train"][0]["question"])  
4 print("Answer: ", raw_datasets["train"][0]["answers"])
```

Context: Beyoncé Giselle Knowles-Carter (/bi:'jɒnseɪ/ bee-YON-say) (born September 4, 1981)
Question: When did Beyonce start becoming popular?
Answer: {'text': ['in the late 1990s'], 'answer_start': [269]}

Module: Reader

❖ Step 4: Create tokenizer

```
{'vincent': 32,  
 'vietnam': 31,  
 'van': 30,  
 'university': 29,  
 'what': 11,  
 '<sep>': 4,  
 '<bos>': 2,  
 'of': 24,  
 'am': 10,  
 'my': 7,  
 'is': 6,  
 'at': 19,  
 '.': 8,  
 'gogh': 22,  
 '<eos>': 3,  
 '<pad>': 1,  
 'computer': 20,  
 'painting': 25,  
 'and': 12,  
 '<unk>': 0,  
 'artist': 18,  
 'favorite': 13,  
 'studying': 15,  
 'i': 5,
```

['I', 'love', 'AIVN']

[5, 23, 17]

```
1 # Load tokenizer để run một số example  
2 tokenizer = AutoTokenizer.from_pretrained(MODEL_NAME)
```

tokenizer_config.json: 100% [██████████] 28.0/28.0 [00:00<00:00, 810B/s]
config.json: 100% [██████████] 483/483 [00:00<00:00, 5.10kB/s]
vocab.txt: 100% [██████████] 232k/232k [00:00<00:00, 973kB/s]
tokenizer.json: 100% [██████████] 466k/466k [00:00<00:00, 1.91MB/s]

Module: Reader

❖ Step 4: Create tokenizer

[CLS]	This	is	the	question	[SEP]	This	is	the	context
with	lots	of	info	##rma	##tion	.	Some	use	##less
.	The	answer	is	here	some	more	words	.	[SEP]

BERT Input format for QA

Module: Reader

❖ Step 4: Create tokenizer

```
1 # Lấy ra 1 example từ tập train
2 context = raw_datasets["train"][0]["context"]
3 question = raw_datasets["train"][0]["question"]
```

```
1 # Sử dụng tokenizer để mã hóa dữ liệu đầu vào
2 inputs = tokenizer(
3     question, # Danh sách các câu hỏi
4     context, # Nội dung liên quan đến câu hỏi
5     max_length=MAX_LENGTH, # Độ dài tối đa cho đầu ra mã hóa
6     truncation="only_second", # Cắt bớt dữ liệu chỉ cho phần thứ hai (context)
7     stride=STRIDE, # Độ dài bước nhảy trong trường hợp dữ liệu dài hơn max_length
8     return_overflowing_tokens=True, # Trả về các tokens vượt quá độ dài tối đa
9     return_offsets_mapping=True, # Trả về bản đồ vị trí của các tokens trong văn bản gốc
10    padding="max_length" # Điền vào dữ liệu để có cùng độ dài max_length
11 )
```

Module: Reader

❖ Step 4: Create tokenizer

```
1 print(f'Question: {question}')  
2 print()  
3 print(f'Context: {context}')
```

Question: When did Beyonce start becoming popular?

Context: Beyoncé Giselle Knowles-Carter (/biː'jɒnseɪ/ bee-YON-say) (born September 4, 1981)

```
1 inputs.keys()
```

```
dict_keys(['input_ids', 'attention mask', 'offset mapping', 'overflow to sample mapping']))
```

1 inputs

Module: Reader

❖ Step 4: Create tokenizer

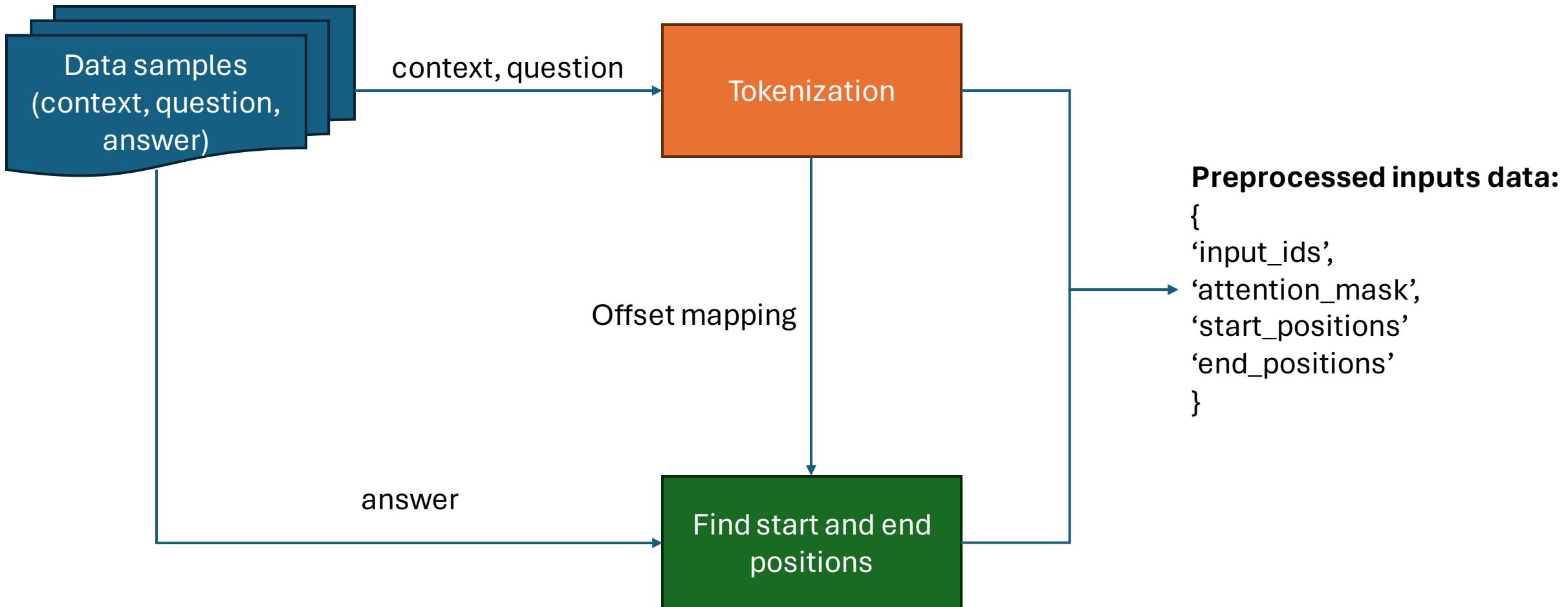
[CLS] question [SEP] context [SEP]

```
1 tokenizer.decode(inputs['input_ids'][0])
```

```
'[CLS] when did beyonce start becoming popular? [SEP] beyonce giselle knowles – carter ( / bi:'jɒnseɪ / bee – yon – say )  
singer, songwriter, record producer and actress. born and raised in houston, texas, she performed in various singing and d  
o fame in the late 1990s as lead singer of r & b girl – group destiny\'s child. managed by her father, mathew knowles, the  
elling girl groups of all time. their hiatus saw the release of beyonce\'s debut album, dangerously in love ( 2003 ), whic  
de, earned five grammy awards and featured the billboard hot 100 number – one singles " crazy in love " and " baby boy ".  
[PAD]  
[AD] [PAD] [PAD] [...'
```

Module: Reader

❖ Step 5: Create preprocessing function



Module: Reader

❖ Step 5: Create preprocessing function

```
1 # Định nghĩa hàm preprocess_training_examples và nhận đối số examples là dữ liệu đào tạo
2 def preprocess_training_examples(examples):
3     # Trích xuất danh sách câu hỏi từ examples và loại bỏ các khoảng trắng dư thừa
4     questions = [q.strip() for q in examples["question"]]
5
6     # Tiến hành mã hóa thông tin đầu vào sử dụng tokenizer
7     inputs = tokenizer(
8         questions,
9         examples["context"],
10        max_length=MAX_LENGTH,
11        truncation="only_second",
12        stride=STRIDE,
13        return_overflowing_tokens=True,
14        return_offsets_mapping=True,
15        padding="max_length",
16    )
17
18     # Trích xuất offset_mapping từ inputs và loại bỏ nó ra khỏi inputs
19     offset_mapping = inputs.pop("offset_mapping")
20
21     # Trích xuất sample_map từ inputs và loại bỏ nó ra khỏi inputs
22     sample_map = inputs.pop("overflow_to_sample_mapping")
23
24     # Trích xuất thông tin về câu trả lời (answers) từ examples
25     answers = examples["answers"]
```

Module: Reader

❖ Step 5: Create preprocessing function

```
# Trích xuất thông tin về câu trả lời (answers) từ examples
answers = examples["answers"]

# Khởi tạo danh sách các vị trí bắt đầu và kết thúc câu trả lời
start_positions = []
end_positions = []

# Duyệt qua danh sách offset_mapping
for i, offset in enumerate(offset_mapping):
    # Xác định index của mẫu (sample) liên quan đến offset hiện tại
    sample_idx = sample_map[i]

    # Trích xuất sequence_ids từ inputs
    sequence_ids = inputs.sequence_ids(i)

    # Xác định vị trí bắt đầu và kết thúc của ngữ cảnh
    idx = 0
    while sequence_ids[idx] != 1:
        idx += 1
    context_start = idx
    while sequence_ids[idx] == 1:
        idx += 1
    context_end = idx - 1

    # Trích xuất thông tin về câu trả lời cho mẫu này
    answer = answers[sample_idx]
```

Finding start and end positions for each sample

Module: Reader

❖ Step 5: Create preprocessing function

```
# Trích xuất thông tin về câu trả lời cho mẫu này
answer = answers[sample_idx]

if len(answer['text']) == 0:
    start_positions.append(0)
    end_positions.append(0)
else:
    # Xác định vị trí bắt đầu và kết thúc của câu trả lời trong ngữ cảnh
    start_char = answer["answer_start"][0]
    end_char = answer["answer_start"][0] + len(answer["text"])[0]

    # Nếu câu trả lời không nằm hoàn toàn trong ngữ cảnh, gán nhãn là (0, 0)
    if offset[context_start][0] > start_char or offset[context_end][1] < end_char:
        start_positions.append(0)
        end_positions.append(0)
    else:
        # Nếu không, gán vị trí bắt đầu và kết thúc dựa trên vị trí của các mã thông tin
        idx = context_start
        while idx <= context_end and offset[idx][0] <= start_char:
            idx += 1
        start_positions.append(idx - 1)

        idx = context_end
        while idx >= context_start and offset[idx][1] >= end_char:
            idx -= 1
        end_positions.append(idx + 1)

# Thêm thông tin vị trí bắt đầu và kết thúc vào inputs
inputs["start_positions"] = start_positions
inputs["end_positions"] = end_positions
```

Finding start and end positions for each sample

Module: Reader

❖ Step 5: Create preprocessing function

```
5 train_dataset = raw_datasets["train"].map(  
6     preprocess_training_examples,  
7     batched=True,  
8     remove_columns=raw_datasets["train"].column_names,  
9 )  
10  
11 # In ra độ dài của tập dữ liệu "train" ban đầu và độ dài của tập dữ liệu đã được xử lý (train_dataset)  
12 len(raw_datasets["train"]), len(train_dataset)
```

Map: 100%  130319/130319 [01:33<00:00, 1327.24 examples/s]
(130319, 131754)

```
1 print(train_dataset[0:10]['start_positions'])  
2 print(train_dataset[0:10]['end_positions'])
```

```
[75, 68, 143, 58, 78, 94, 134, 101, 77, 84]  
[78, 70, 143, 60, 79, 97, 136, 102, 78, 85]
```

Module: Reader

❖ Step 5: Create preprocessing function

```
1 def preprocess_validation_examples(examples):          17 # Lấy ánh xạ để ánh xạ lại ví dụ tham chiếu cho từng dòng trong inputs
2     # Chuẩn bị danh sách câu hỏi bằng cách loại bỏ khoảng trống
3     questions = [q.strip() for q in examples["question"]]
4
5     # Sử dụng tokenizer để mã hóa các câu hỏi và văn bản liên quan
6     inputs = tokenizer(
7         questions,
8         examples["context"],
9         max_length=MAX_LENGTH,
10        truncation="only_second",
11        stride=STRIDE,
12        return_overflowing_tokens=True,
13        return_offsets_mapping=True,
14        padding="max_length",
15    )
16
17     # Lấy ánh xạ để ánh xạ lại ví dụ tham chiếu cho từng dòng trong inputs
18     sample_map = inputs.pop("overflow_to_sample_mapping")
19     example_ids = []
20
21     # Xác định ví dụ tham chiếu cho mỗi dòng đầu vào và điều chỉnh ánh xạ offset
22     for i in range(len(inputs["input_ids"])):
23         sample_idx = sample_map[i]
24         example_ids.append(examples["id"][sample_idx])
25
26         sequence_ids = inputs.sequence_ids(i)
27         offset = inputs["offset_mapping"][i]
28
29         # Loại bỏ các offset không phù hợp với sequence_ids
30         inputs["offset_mapping"][i] = [
31             o if sequence_ids[k] == 1 else None for k, o in enumerate(offset)
32         ]
33
34     # Thêm thông tin ví dụ tham chiếu vào đầu vào
35     inputs["example_id"] = example_ids
36
37     return inputs
```

Module: Reader

❖ Step 5: Create preprocessing function

```
2 validation_dataset = raw_datasets["validation"].map(  
3     preprocess_validation_examples, # Gọi hàm preprocess_validation_examples  
4     batched=True, # Xử lý dữ liệu theo từng batch.  
5     remove_columns=raw_datasets["validation"].column_names, # Loại bỏ các cột  
6 )  
7  
8 # In ra độ dài của raw_datasets["validation"] và validation_dataset để so sánh  
9 len(raw_datasets["validation"]), len(validation_dataset)
```

Map: 100%

11873/11873 [00:14<00:00, 938.54 examples/s]

(11873, 12134)

Module: Reader

❖ Step 6: Training

```
1 # Load model
2 model = AutoModelForQuestionAnswering.from_pretrained(MODEL_NAME)

model.safetensors: 100% [268M/268M [00:01<00:00, 170MB/s]

Some weights of DistilBertForQuestionAnswering were not initialized from the model checkpoint
You should probably TRAIN this model on a down-stream task to be able to use it for prediction
```

```
1 # Tạo đối tượng args là các tham số cho quá trình huấn luyện
2 args = TrainingArguments(
3     output_dir="distilbert-finetuned-squadv2", # Thư mục lưu trữ kết quả huấn luyện
4     evaluation_strategy="no", # Chế độ đánh giá không tự động sau mỗi epoch
5     save_strategy="epoch", # Lưu checkpoint sau mỗi epoch
6     learning_rate=2e-5, # Tốc độ học
7     num_train_epochs=3, # Số epoch huấn luyện
8     weight_decay=0.01, # Giảm trọng lượng mô hình để tránh overfitting
9     fp16=True, # Sử dụng kiểu dữ liệu half-precision (16-bit) để tối ưu hóa tài nguyên
10    push_to_hub=True, # Đẩy kết quả huấn luyện lên một nơi chia sẻ trực tuyến (Hub)
11 )
```

```
1 # Khởi tạo một đối tượng Trainer để huấn luyện mô hình
2 trainer = Trainer(
3     model=model, # Sử dụng mô hình đã tạo trước đó
4     args=args, # Các tham số và cấu hình huấn luyện
5     train_dataset=train_dataset, # Sử dụng tập dữ liệu huấn luyện
6     eval_dataset=validation_dataset, # Sử dụng tập dữ liệu đánh giá
7     tokenizer=tokenizer, # Sử dụng tokenizer để xử lý văn bản
8 )
9
10 # Bắt đầu quá trình huấn luyện
11 trainer.train()
```

You're using a DistilBertTokenizerFast tokenizer.
[49354/49410 1:3]

Step	Training Loss
500	3.102900
1000	2.268800
1500	1.971300
2000	1.810200
2500	1.702800
3000	1.625000
3500	1.574800
4000	1.566400
4500	1.495000
5000	1.480800
5500	1.438300
6000	1.402200
6500	1.418000
7000	1.426900
7500	1.366800
8000	1.351200
8500	1.318100
9000	1.307100
9500	1.308100
10000	1.284500

Module: Reader

❖ Step 7: Test

```
predicted_answers = []
for example in tqdm(examples):
    example_id = example['id']
    context = example['context']
    answers = []

    # Lặp qua tất cả các đặc trưng liên quan đến ví dụ đó
    for feature_index in example_to_features[example_id]:
        start_logit = start_logits[feature_index]
        end_logit = end_logits[feature_index]
        offsets = features[feature_index]['offset_mapping']

        # Lấy các chỉ số có giá trị lớn nhất cho start và end logits
        start_indexes = np.argsort(start_logit)[-1 : -N_BEST - 1 : -1].tolist()
        end_indexes = np.argsort(end_logit)[-1 : -N_BEST - 1 : -1].tolist()
        for start_index in start_indexes:
            for end_index in end_indexes:
                # Bỏ qua các câu trả lời không hoàn toàn nằm trong ngữ cảnh
                if offsets[start_index] is None or offsets[end_index] is None:
                    continue
                # Bỏ qua các câu trả lời có độ dài > max_answer_length
                if end_index - start_index + 1 > MAX_ANS_LENGTH:
                    continue

                # Tạo một câu trả lời mới
                answer = {
                    'text': context[offsets[start_index][0] : offsets[end_index][1]],
                    'logit_score': start_logit[start_index] + end_logit[end_index],
                }
                answers.append(answer)

# Chọn câu trả lời có điểm số tốt nhất
if len(answers) > 0:
    best_answer = max(answers, key=lambda x: x['logit_score'])
    answer_dict = {
        'id': example_id,
        'prediction_text': best_answer['text'],
        'no_answer_probability': 1 - best_answer['logit_score']
    }
else:
    answer_dict = {
        'id': example_id,
        'prediction_text': '',
        'no_answer_probability': 1.0
    }
predicted_answers.append(answer_dict)
```

```
# Chọn câu trả lời có điểm số tốt nhất
if len(answers) > 0:
    best_answer = max(answers, key=lambda x: x['logit_score'])
    answer_dict = {
        'id': example_id,
        'prediction_text': best_answer['text'],
        'no_answer_probability': 1 - best_answer['logit_score']
    }
else:
    answer_dict = {
        'id': example_id,
        'prediction_text': '',
        'no_answer_probability': 1.0
    }
predicted_answers.append(answer_dict)
```

100%

```
{'exact': 47.452202476206516,  
 'f1': 51.28265600122848,  
 'total': 11873,  
 'HasAns_exact': 74.78070175438596,  
 'HasAns_f1': 82.45259357331055,  
 'HasAns_total': 5928,  
 'NoAns_exact': 20.201850294365013,  
 'NoAns_f1': 20.201850294365013,  
 'NoAns_total': 5945,  
 'best_exact': 64.46559420533984,  
 'best_exact_thresh': -11.35546875,  
 'best_f1': 66.16170764530801,  
 'best_f1_thresh': -9.61328125}
```

Module: Reader

❖ Step 7: Test

```
1 # Use a pipeline as a high-level helper
2 from transformers import pipeline
3
4 PIPELINE_NAME = 'question-answering'
5 MODEL_NAME = 'thangduong0509/distilbert-finetuned-squadv2'
6 pipe = pipeline(PIPELINE_NAME, model=MODEL_NAME)
```

```
1 INPUT_QUESTION = 'What is my name?'
2 INPUT_CONTEXT = 'My name is AI Vietnam and I live in Vietnam.'
3 pipe(question=INPUT_QUESTION, context=INPUT_CONTEXT)
```

```
{'score': 0.97179114818573, 'start': 11, 'end': 21, 'answer': 'AI Vietnam'}
```

Module: Reader

❖ Step 7: Test

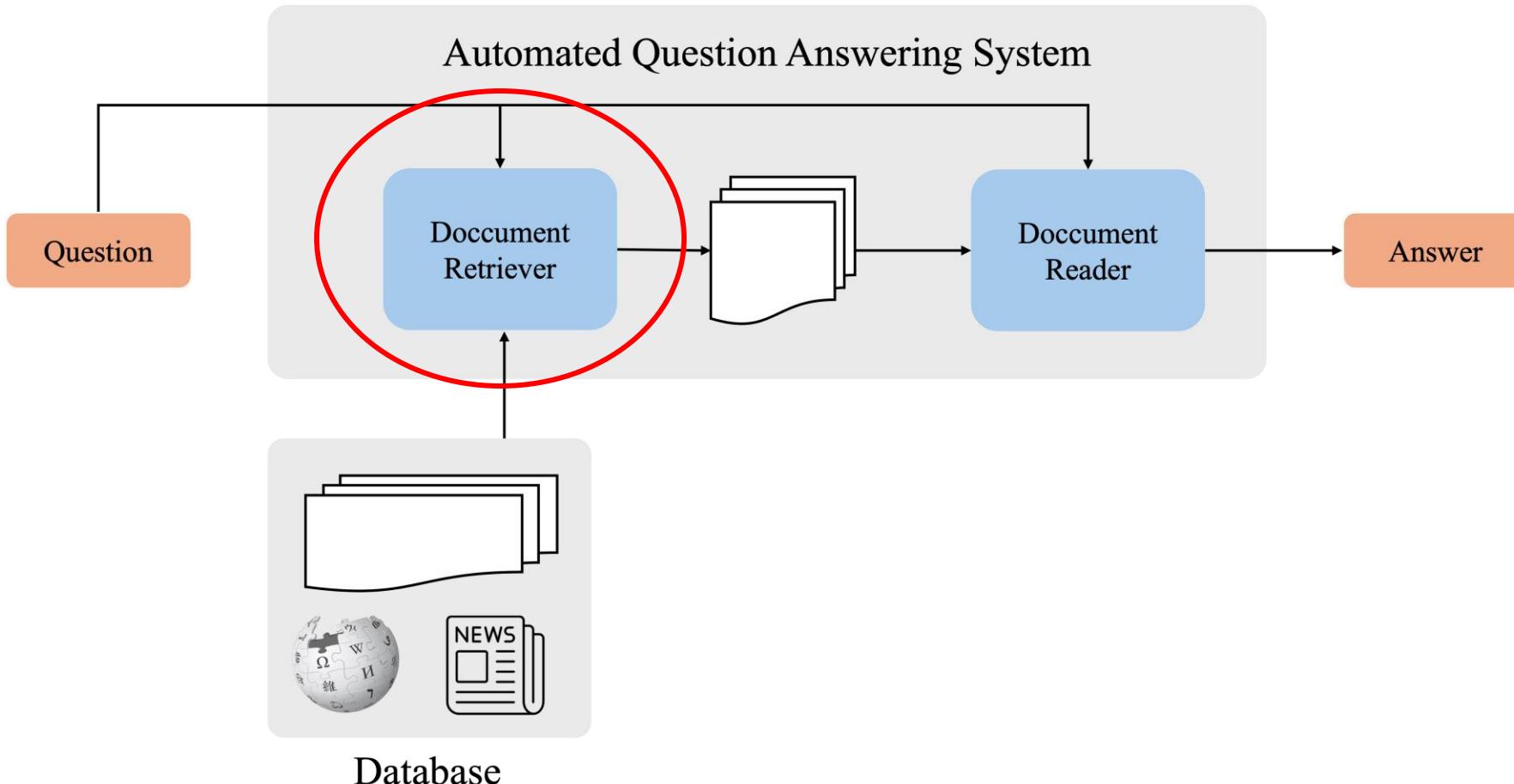
```
1 INPUT_QUESTION = 'where is the highest mountain in solar system located'  
2 INPUT_CONTEXT = 'The highest mountain and volcano in the Solar System is on the planet Mars.  
3 pipe(question=INPUT_QUESTION, context=INPUT_CONTEXT)
```

```
{'score': 0.852756679058075, 'start': 70, 'end': 74, 'answer': 'Mars'}
```

Module: Retriever

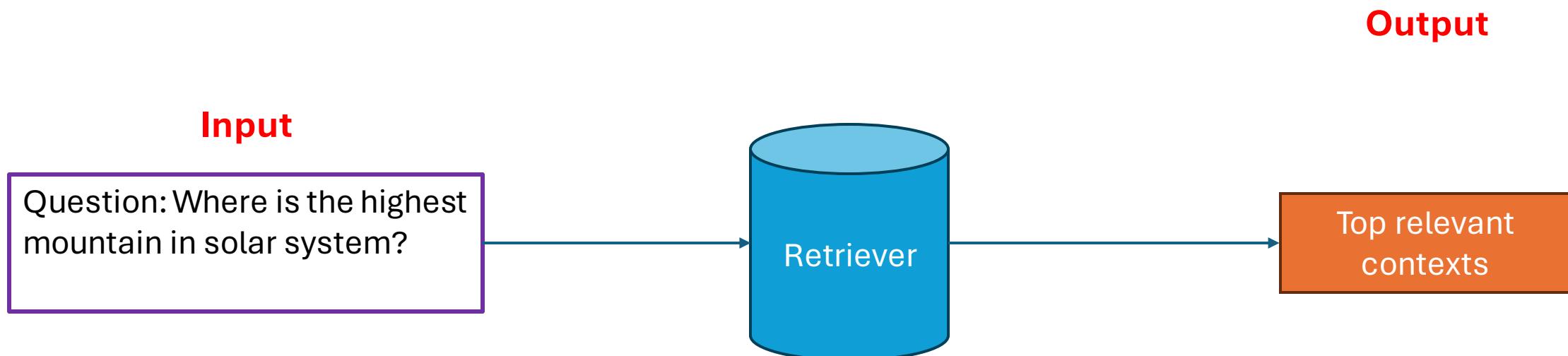
Module: Retriever

❖ E2E QA Pipeline



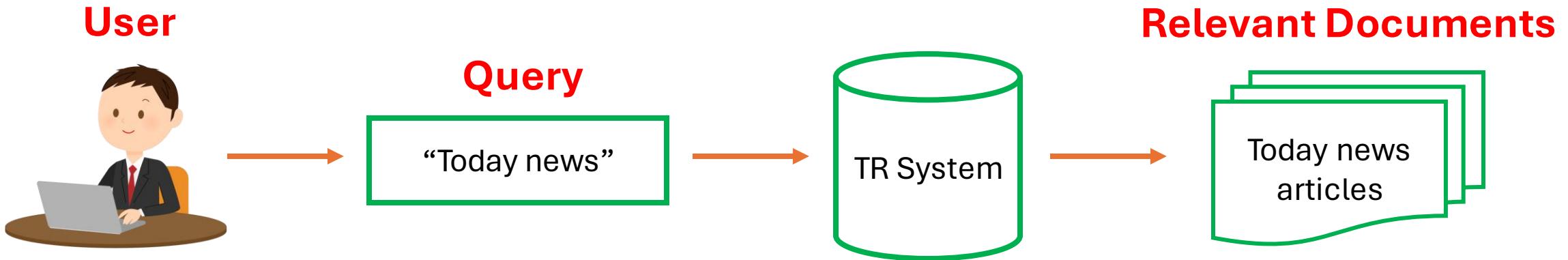
Module: Retriever

❖ Retriever I/O



Module: Retriever

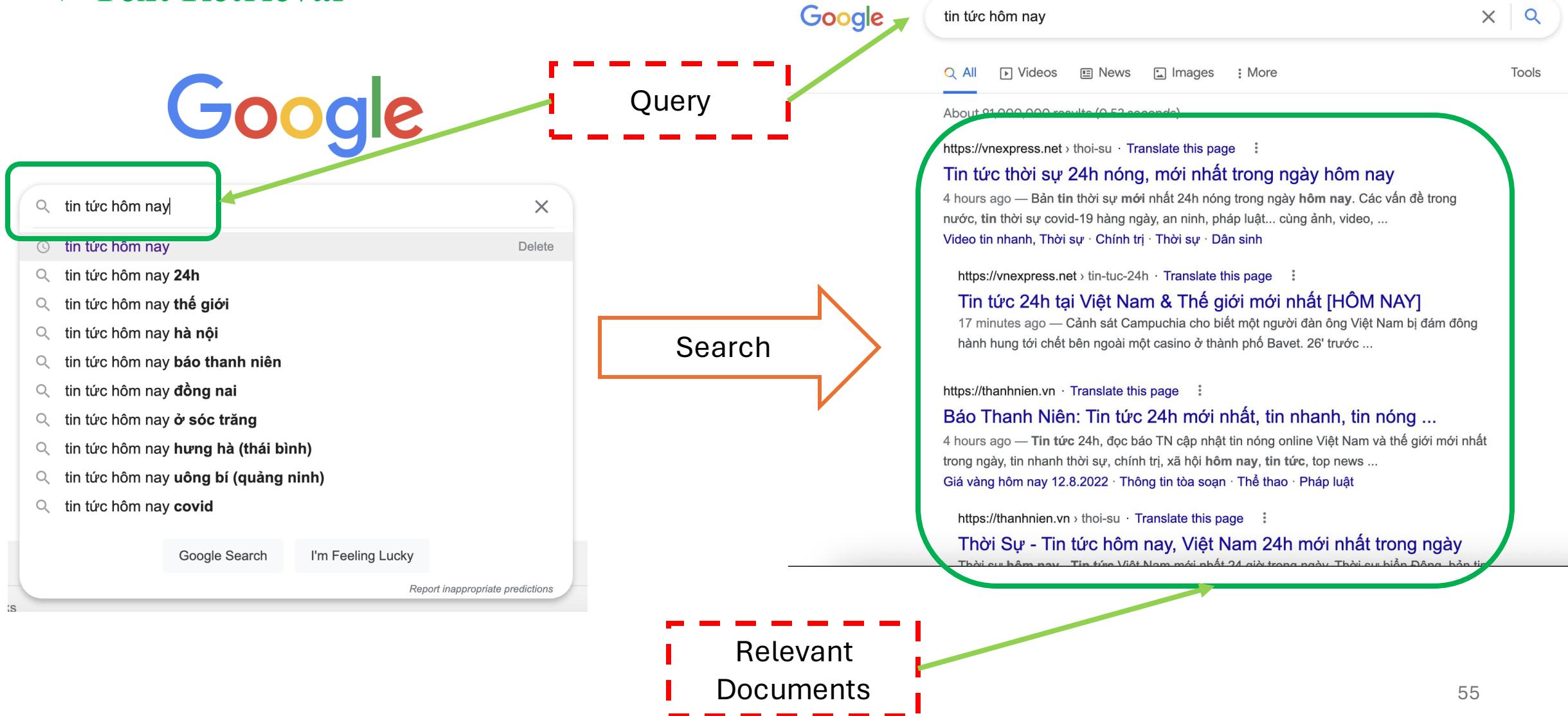
❖ Text Retrieval



Text Retrieval (TR) (also called as Document Retrieval)¹: A branch of Information Retrieval (IR) where the system matching of some stated user search query against a set of texts.

Module: Retriever

❖ Text Retrieval

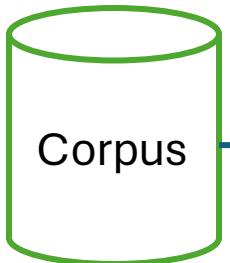


Module: Retriever

❖ Text Retriever Pipeline

'what is the official language
in Fiji'

Query (Question)



SQuAD_V2 Corpus

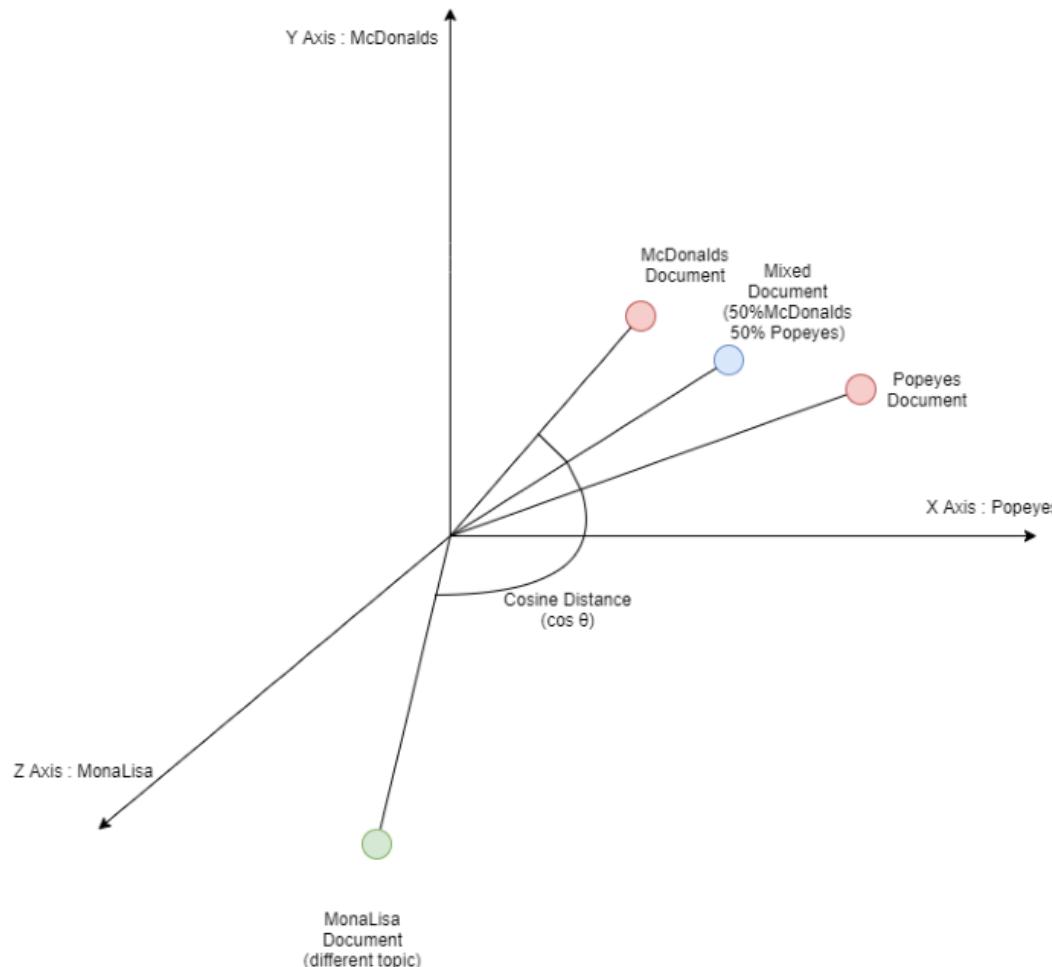
Retriever

- Query: what is the official language in Fiji
==== Relevant docs ===
Top 1; Score: 0.6556
The official languages in Fiji are Fijian and English. A dialect of Hindustani is als
- Top 2; Score: 0.6556
The official languages in Fiji are Fijian and English. A dialect of Hindustani is als
- Top 3; Score: 0.5715
The official languages. Fiji's 1997 Constitution established Fijian as one of the off
- Top 4; Score: 0.5604
Of all the languages of Russia, Russian is the only official language. There are 35 c
- Top 5; Score: 0.5592
Finnish is one of two official languages of Finland (the other being Swedish, spoken
- Top 6; Score: 0.5307
Liberia is a multilingual country where more than thirty languages are spoken. English
- Top 7; Score: 0.5164
There are over 120 ethnic groups, each with its own language or dialect. Indigenous T
- Top 8; Score: 0.5094
the two official languages every country in the world has official languages that are
- Top 9; Score: 0.4811
Fiji Language. The Fijian language spoken in Fiji is a type of Austronesian langauge
- Top 10; Score: 0.4547
Upon Uganda's independence in 1962, English was maintained as the official language,

Relevant Questions

Module: Retriever

❖ Idea

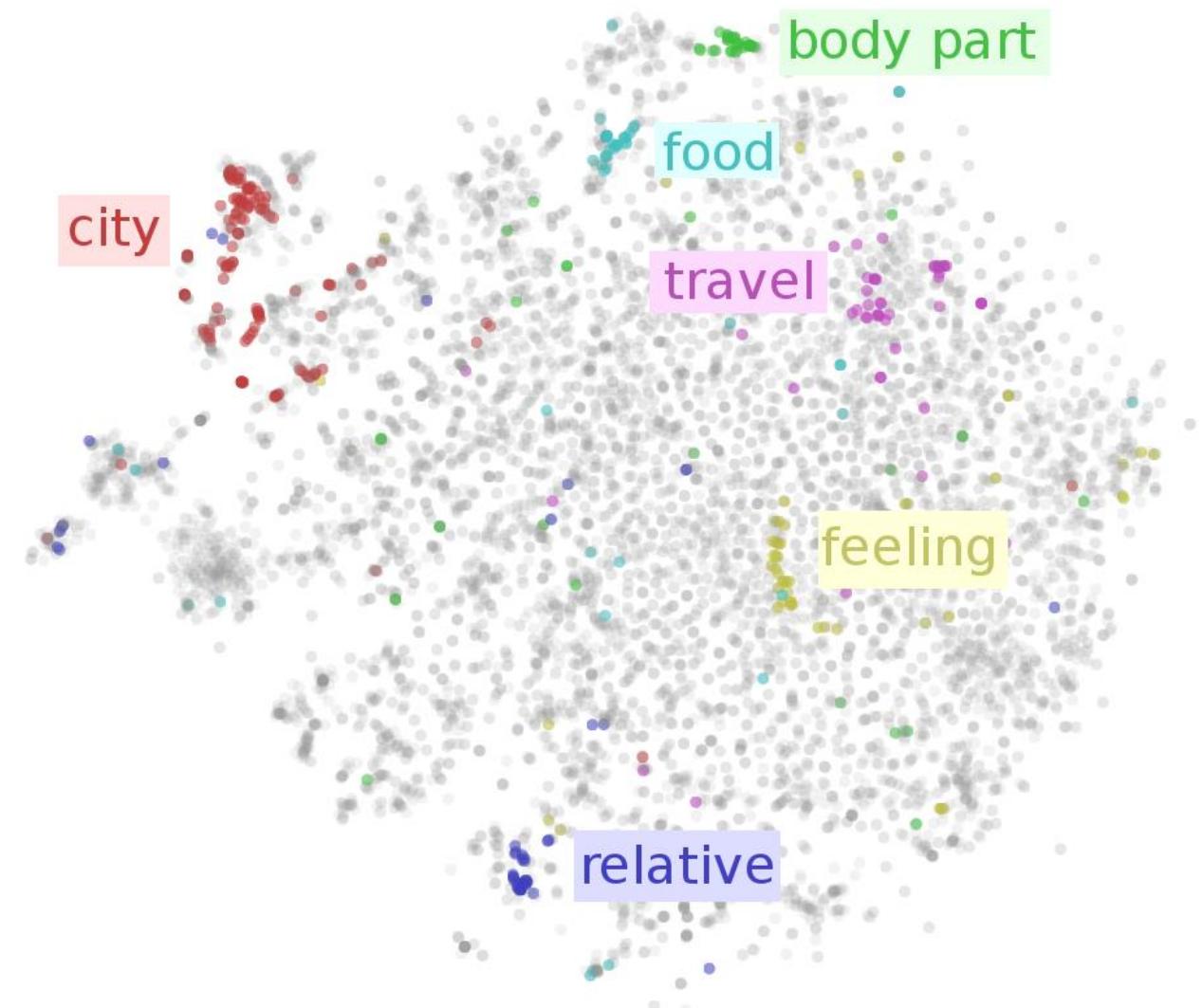
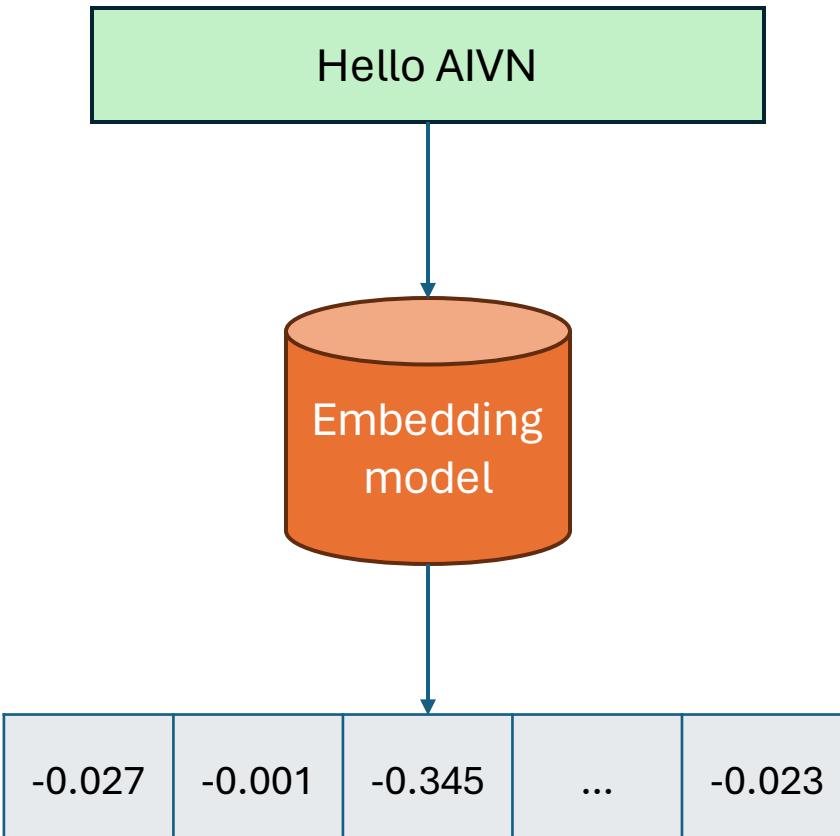


Given two vectors, we can calculate the similarity between them by using a similarity formula. E.g: L2 Distance, Cosine Similarity...

$$\text{similarity}(A, B) = \frac{A \cdot B}{\|A\| \times \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n A_i^2} \times \sqrt{\sum_{i=1}^n B_i^2}}$$

Module: Retriever

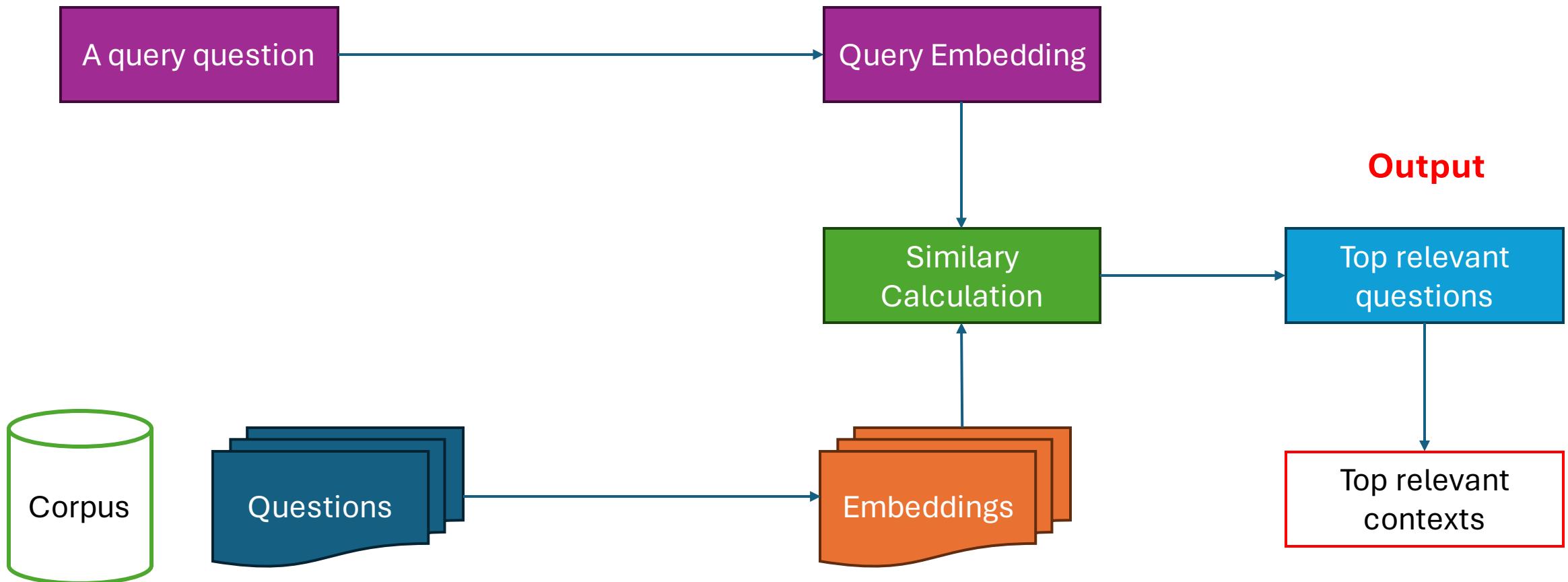
❖ Idea



Module: Retriever

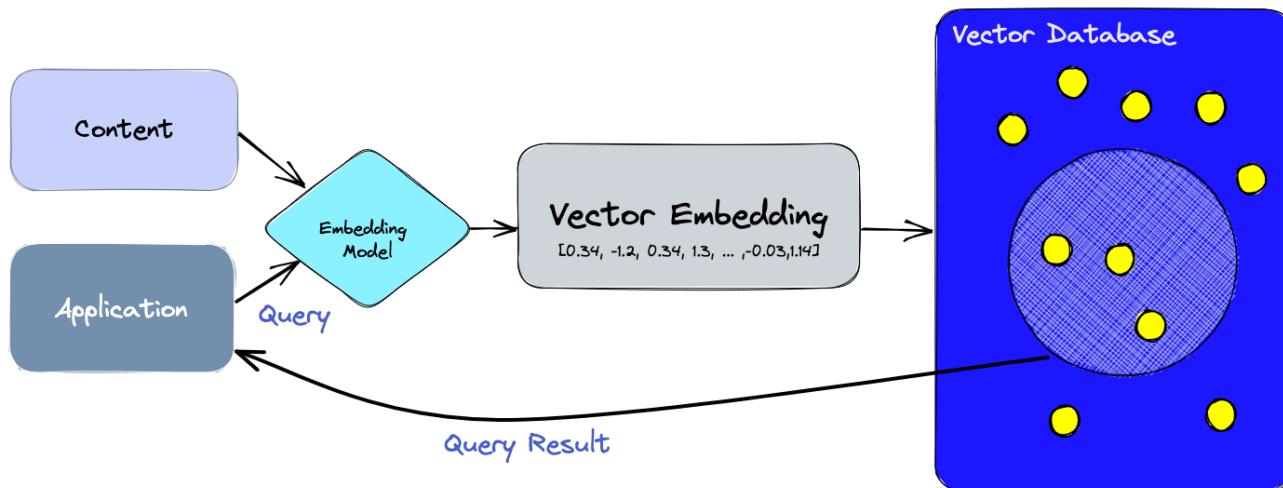
❖ Idea

Input



Module: Retriever

❖ Vector Database



Vector database: A database that can store vectors (fixed-length lists of numbers) along with other data items.

Module: Retriever

❖ Vector Database

Choosing a Vector Database

Vector Databases	Vector Libraries
 Milvus	 Pinecone
 vespa	 chroma
 LanceDB	 Weaviate
 marqo	 drant
	 Vald
Vector-Capable NoSQL Databases	Vector-Capable SQL Databases
 mongoDB	 neo4j
 ROCKSET	 redis
 cassandra	 Timescale
 DataStax Astra	 SingleStore
	 kinetica
	 ClickHouse
Text Search Databases	
	 elasticsearch
	 Solr
	 OpenSearch
	 LUCENE

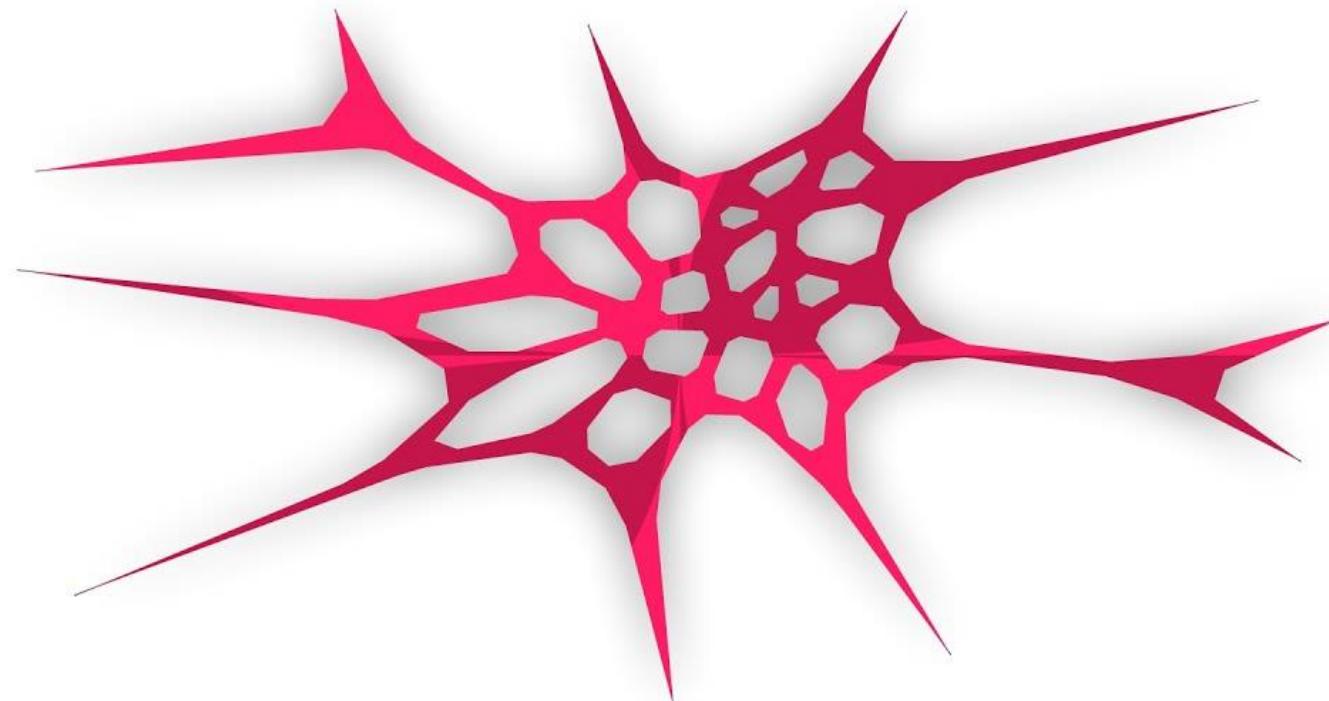
graft

Module: Retriever

❖ Faiss

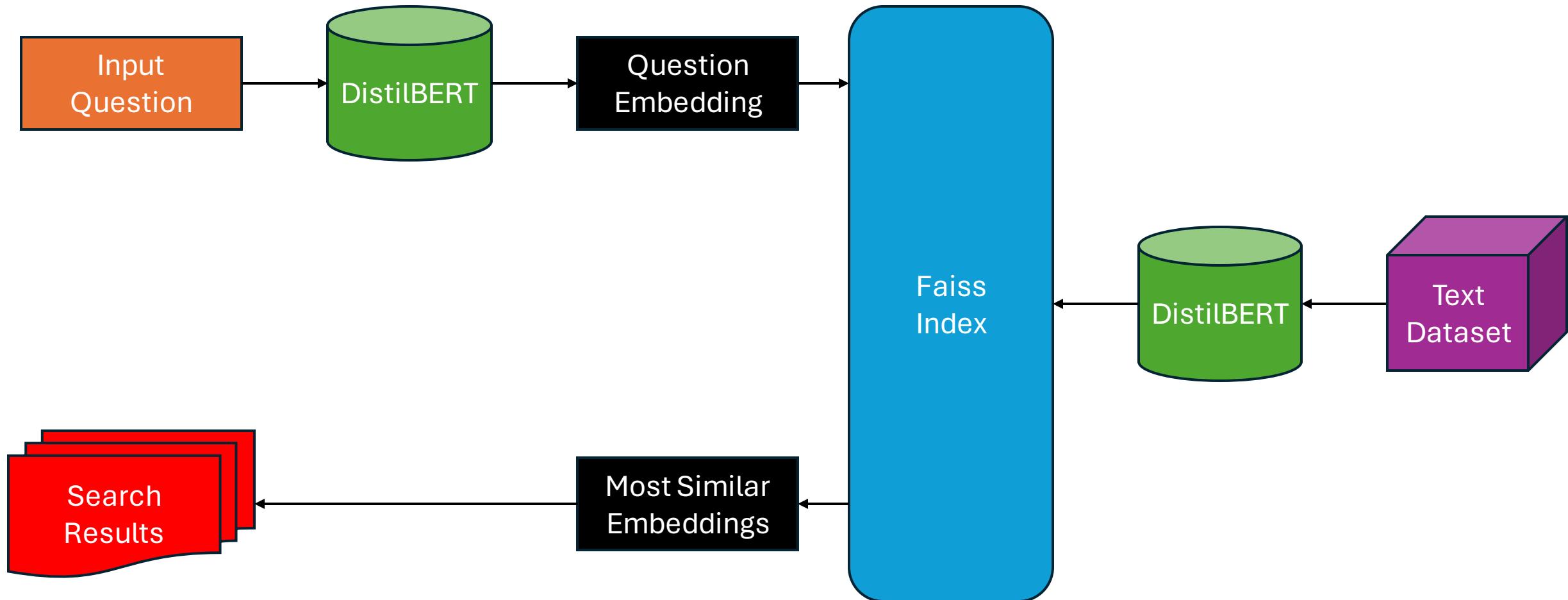
FAISS

Scalable Search With Facebook AI



Module: Retriever

❖ Faiss Pipeline



Module: Retriever

❖ Step 1: Install and import libraries

```
1 !pip install -qq transformers[sentencepiece]==4.35.2 datasets==2.16.1 evaluate==0.4.1
```

```
████████████████████ 507.1/507.1 kB 7.0 MB/s eta 0:00:00
████████████████████ 84.1/84.1 kB 6.6 MB/s eta 0:00:00
████████████████████ 115.3/115.3 kB 10.6 MB/s eta 0:00:00
████████████████████ 134.8/134.8 kB 8.1 MB/s eta 0:00:00
████████████████████ 134.8/134.8 kB 10.6 MB/s eta 0:00:00
```

```
1 !sudo apt-get install libomp-dev
2 !pip install -qq faiss-gpu
```

Module: Retriever

❖ Step 1: Install and import libraries

```
1 import numpy as np
2 import collections
3 import torch
4 import faiss
5 import evaluate
6
7 from datasets import load_dataset
8 from transformers import AutoTokenizer, AutoModel
9 from transformers import AutoModelForQuestionAnswering
10 from transformers import TrainingArguments
11 from transformers import Trainer
12 from tqdm.auto import tqdm
13
14 device = torch.device("cuda") if torch.cuda.is_available() else torch.device("cpu")
```

Module: Retriever

❖ Step 2: Download dataset

```
1 DATASET_NAME = 'squad_v2'  
2 raw_datasets = load_dataset(DATASET_NAME, split='train')  
3 raw_datasets  
  
/usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_token.py:88: UserWarning:  
The secret `HF_TOKEN` does not exist in your Colab secrets.  
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens),  
You will be able to reuse this secret in all of your notebooks.  
Please note that authentication is recommended but still optional to access public models or datasets.  
    warnings.warn(  
Downloading readme: 100%  8.18k/8.18k [00:00<00:00, 376kB/s]  
Downloading data: 100%  16.4M/16.4M [00:05<00:00, 2.81MB/s]  
Downloading data: 100%  1.35M/1.35M [00:04<00:00, 309kB/s]  
Generating train split: 100%  130319/130319 [00:00<00:00, 265926.85 examples/s]  
Generating validation split: 100%  11873/11873 [00:00<00:00, 99856.77 examples/s]  
Dataset({  
    features: ['id', 'title', 'context', 'question', 'answers'],  
    num_rows: 130319  
})
```

Module: Retriever

❖ Step 2: Download dataset

```
1 raw_datasets = raw_datasets.filter(  
2     lambda x: len(x['answers']['text']) > 0  
3 )  
4 raw_datasets
```

Filter: 100%  130319/130319 [00:04<00:00, 30967.75 examples/s]

```
Dataset({  
    features: ['id', 'title', 'context', 'question', 'answers'],  
    num_rows: 86821  
})
```

context string · lengths	question string · lengths
 151 3.71k	 1 25.7k
Beyoncé Giselle Knowles-Carter (/bi...	When did Beyonce start becoming...
Beyoncé Giselle Knowles-Carter (/bi...	What areas did Beyonce compete in...
Beyoncé Giselle Knowles-Carter (/bi...	When did Beyonce leave Destiny's...
Beyoncé Giselle Knowles-Carter (/bi...	In what city and state did Beyonce...
Beyoncé Giselle Knowles-Carter (/bi...	In which decade did Beyonce become...
Beyoncé Giselle Knowles-Carter (/bi...	In what R&B group was she the lead...

Module: Retriever

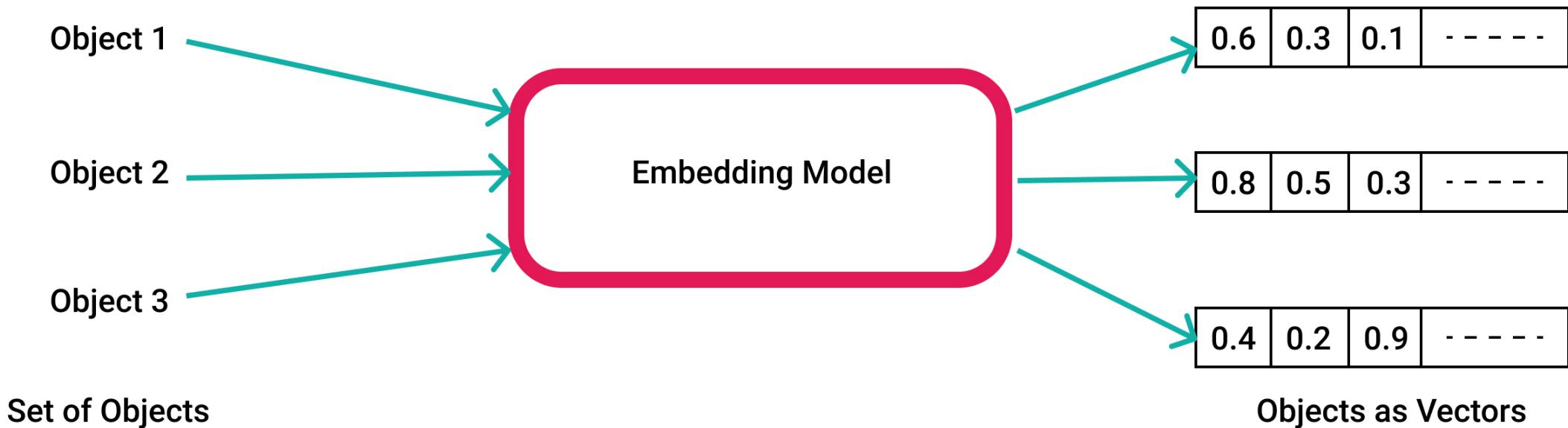
❖ Step 3: Initialize model

```
1 MODEL_NAME = "distilbert-base-uncased"  
2 tokenizer = AutoTokenizer.from_pretrained(MODEL_NAME)  
3 model = AutoModel.from_pretrained(MODEL_NAME).to(device)
```



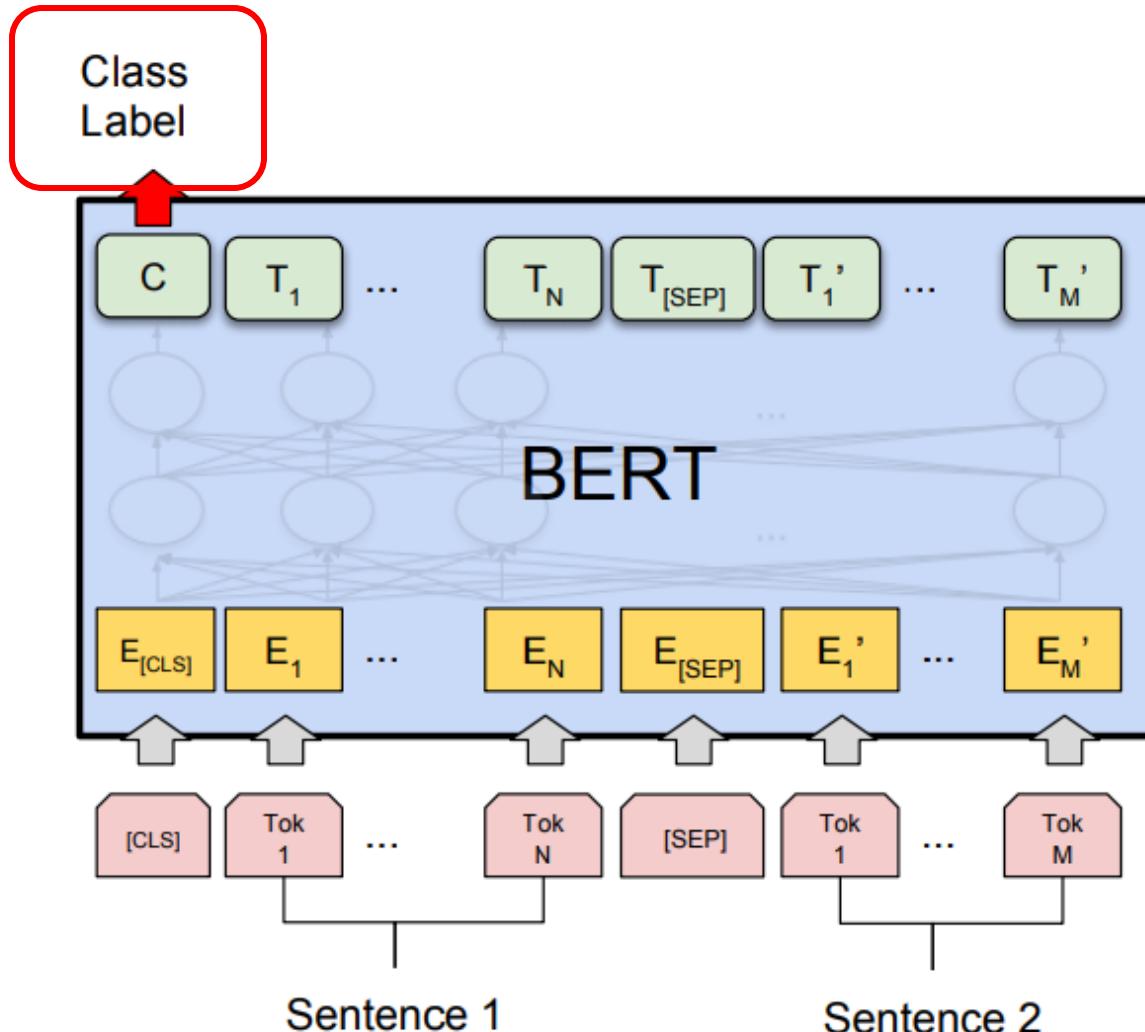
Module: Retriever

❖ Step 4: Create get embedding function



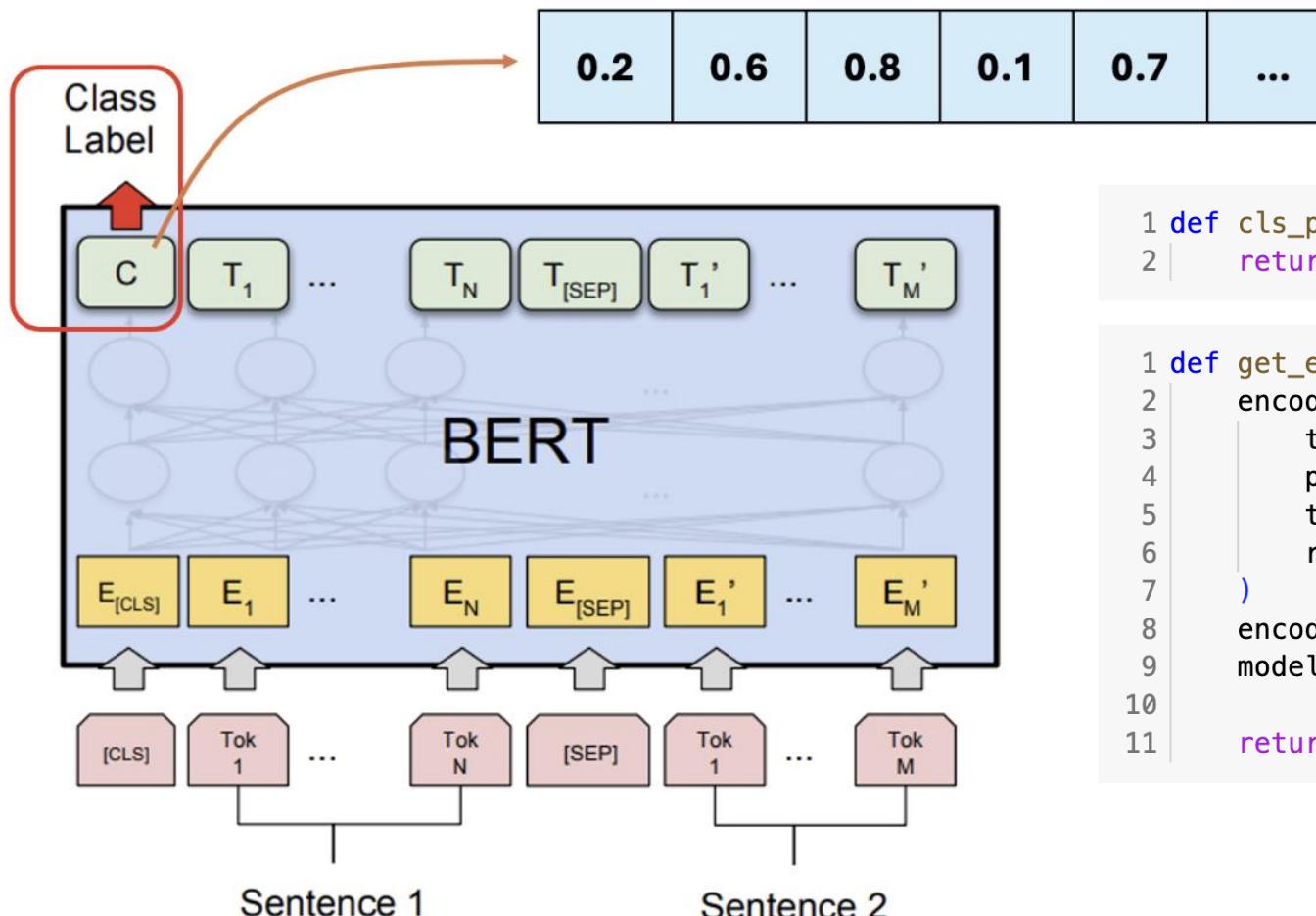
Module: Retriever

❖ Step 4: Create get embedding function



Module: Retriever

❖ Step 4: Create get embedding function

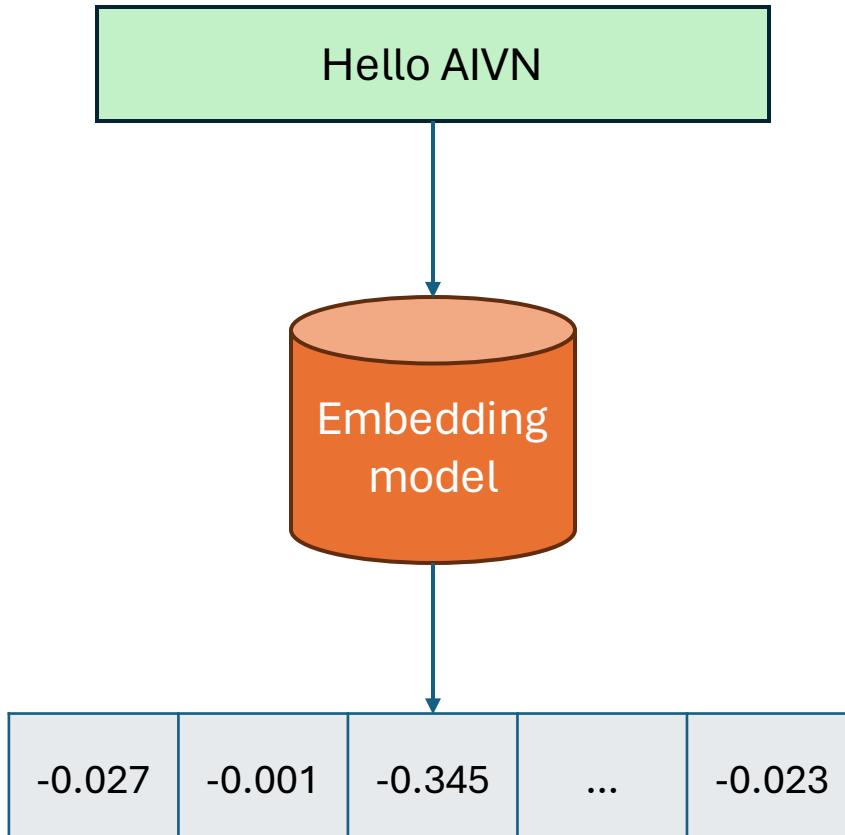


```
1 def cls_pooling(model_output):
2     return model_output.last_hidden_state[:, 0]
3
4
5
6
7
8
9
10
11
```

```
1 def get_embeddings(text_list):
2     encoded_input = tokenizer(
3         text_list,
4         padding=True,
5         truncation=True,
6         return_tensors='pt'
7     )
8     encoded_input = {k: v.to(device) for k, v in encoded_input.items()}
9     model_output = model(**encoded_input)
10
11     return cls_pooling(model_output)
```

Module: Retriever

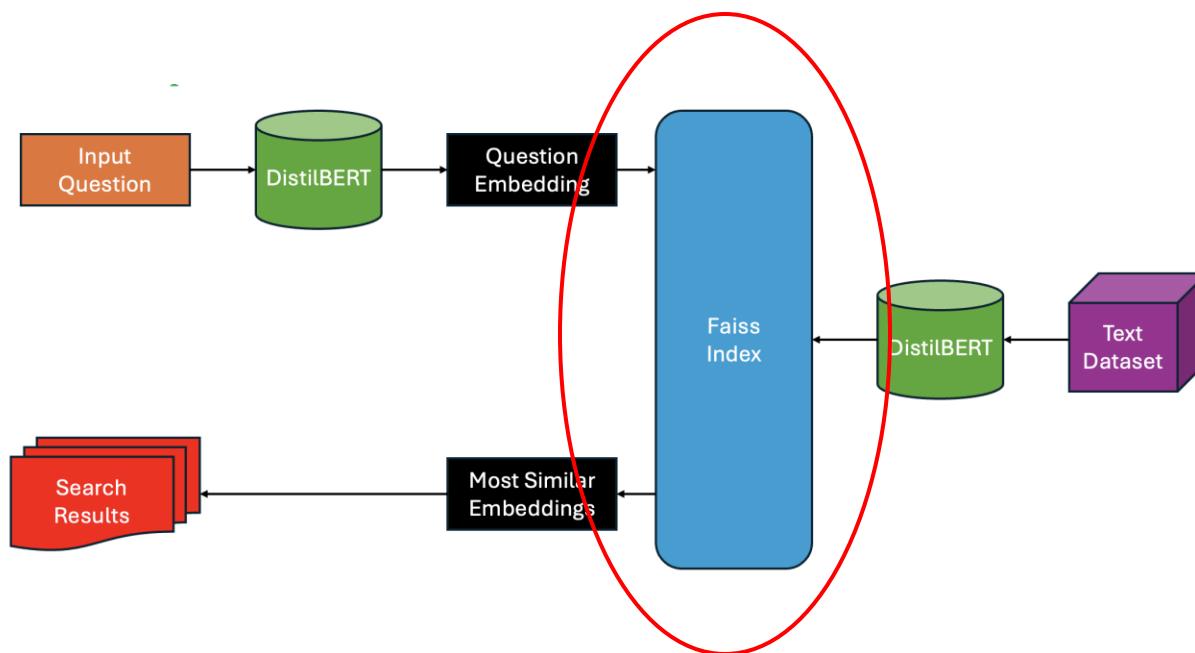
❖ Step 4: Create get embedding function



```
1 # Test functionality
2 embedding = get_embeddings(raw_datasets['question'][0])
3 embedding.shape
torch.Size([1, 768])
```

Module: Retriever

❖ Step 4: Create get embedding function



```
1 # Convert to numpy array (required for HF Datasets)
2 EMBEDDING_COLUMN = 'question_embedding'
3 embeddings_dataset = raw_datasets.map(
4     lambda x: {EMBEDDING_COLUMN: get_embeddings(x['question']).detach().cpu().numpy()[0]}
5 )
```

Map: 100% [86821/86821 [09:42<00:00, 142.06 examples/s]]

```
1 embeddings_dataset.add_faiss_index(column=EMBEDDING_COLUMN)
```

100% [87/87 [00:00<00:00, 109.40it/s]]

```
Dataset({
    features: ['id', 'title', 'context', 'question', 'answers', 'question_embedding'],
    num_rows: 86821
})
```

Module: Retriever

❖ Step 5: Search

```
1 question = 'When did Beyonce start becoming popular?'
2
3 input_quest_embedding = get_embeddings([question]).cpu().detach().numpy()
4 input_quest_embedding.shape
```

(1, 768)

```
1 TOP_K = 5
2 scores, samples = embeddings_dataset.get_nearest_examples(
3     EMBEDDING_COLUMN, input_quest_embedding, k=TOP_K
4 )
```

- | Top 1 Score: 0.0
Question: When did Beyonce start becoming popular?
Context: Beyoncé Giselle Knowles-Carter (/bi:'jɒnseɪ/ bee-YON-say) (bo
Answer: {'text': ['in the late 1990s'], 'answer_start': [269]})
- Top 2 Score: 2.6135313510894775
Question: When did Beyoncé rise to fame?
Context: Beyoncé Giselle Knowles-Carter (/bi:'jɒnseɪ/ bee-YON-say) (bo
Answer: {'text': ['late 1990s'], 'answer_start': [276]})
- Top 3 Score: 4.859482288360596
Question: When did Beyoncé release Formation?
Context: On February 6, 2016, one day before her performance at the Su
Answer: {'text': ['February 6, 2016'], 'answer_start': [3]})
- Top 4 Score: 5.054221153259277
Question: In which decade did Beyonce become famous?
Context: Beyoncé Giselle Knowles-Carter (/bi:'jɒnseɪ/ bee-YON-say) (bo
Answer: {'text': ['late 1990s'], 'answer_start': [276]})
- Top 5 Score: 5.170375347137451
Question: When did Beyonce begin her deals with name brands?
Context: The release of a video-game Starpower: Beyoncé was cancelled
Answer: {'text': ['since the age of 18'], 'answer_start': [433]})

Module: Retriever

❖ Step 5: Search with QA

```
1 question = 'When did Beyonce start becoming popular?'
2
3 input_quest_embedding = get_embeddings([question]).cpu().detach().numpy()
4 input_quest_embedding.shape

(1, 768)

1 TOP_K = 3
2 for idx, input_question in enumerate(embeddings_dataset['question'][200:210]):
3     input_quest_embedding = get_embeddings([input_question]).cpu().detach().numpy()
4     scores, samples = embeddings_dataset.get_nearest_examples(
5         EMBEDDING_COLUMN, input_quest_embedding, k=TOP_K
6     )
7     print(f'Question {idx + 1}: {input_question}')
8     for jdx, score in enumerate(scores):
9         print(f'Top {jdx + 1}\tScore: {score}')
10        question = samples['question'][jdx]
11        context = samples['context'][jdx]
12        answer = pipe(
13            question=question,
14            context=context
15        )
16        print(f'Context: {context}')
17        print(f'Answer: {answer}')
18        print()
19 print()
```

Question 1: How many awards was Beyonce nominated for at the 52nd Grammy Awards?
Top 1 Score: 0.0
Context: At the 52nd Annual Grammy Awards, Beyoncé received ten nominations, includ
Answer: {'score': 0.684537410736084, 'start': 51, 'end': 54, 'answer': 'ten'}

Top 2 Score: 2.3847267627716064
Context: At the 57th Annual Grammy Awards in February 2015, Beyoncé was nominated f
Answer: {'score': 0.3522382378578186, 'start': 77, 'end': 80, 'answer': 'six'}

Top 3 Score: 4.536745548248291
Context: Beyoncé's first solo recording was a feature on Jay Z's "'03 Bonnie & Clyd
Answer: {'score': 0.8402150273323059, 'start': 747, 'end': 751, 'answer': 'five'}

Question 2: Beyonce tied with which artist for most nominations by a female artist?
Top 1 Score: 0.0
Context: At the 52nd Annual Grammy Awards, Beyoncé received ten nominations, including Alb
Answer: {'score': 0.9431203603744507, 'start': 242, 'end': 253, 'answer': 'Laure Hill'}

Top 2 Score: 8.78997802734375
Context: At the 52nd Annual Grammy Awards, Beyoncé received ten nominations, including Alb
Answer: {'score': 0.6438630819320679, 'start': 517, 'end': 529, 'answer': 'Maria Carey'}

Top 3 Score: 11.078895568847656
Context: At the 52nd Annual Grammy Awards, Beyoncé received ten nominations, including Alb
Answer: {'score': 0.9386846423149109, 'start': 242, 'end': 253, 'answer': 'Laure Hill'}

Question

