

**COMP5046**

# *Natural Language Processing*

*Lecture 10: Attention and Question Answering  
(Reading Comprehension)*

*Dr. Caren Han*

*Semester 1, 2022*

*School of Computer Science,  
University of Sydney*



# 0 LECTURE PLAN

## Lecture 10: Attention and Question Answering (Reading Comprehension)

1. Question Answering
2. Knowledge-based Question Answering
3. IR-based Question Answering (Reading Comprehension)
4. Attention
5. Reading Comprehension with Attention
6. Visual Question Answering

# 1 Question Answering

## Question Answering

Major Aim \*  
communicate between human & machine.

**Question answering (QA)** is a computer science discipline within the fields of information retrieval and natural language processing (NLP), which is concerned with building systems that automatically answer questions posed by humans in a natural language.

### Different types of questions:

**General questions**, with Yes/No answers

- e.g. Are you a student?

**Wh- Questions**, start with: who, what, where, when, why, how, how many

- e.g. When did you get to this lecture?
- e.g. What is the weather like in London?



# 1 Question Answering

## Question Answering

**Question answering (QA)** is a computer science discipline within the fields of information retrieval and natural language processing (NLP), which is concerned with building systems that automatically answer questions posed by humans in a natural language.

### *Different types of questions:*

**Choice Questions**, where you have some options inside the question

**Factoid questions**, where the complete answer can be found inside a text. The answer to such questions consist of one or several words that go one after another



# 1 Question Answering

## Question

*Three Questions for building a QA System*

- What do the answers look like? Y / N ? Single word ?
- Where can I get the answers from? find source
- What does my training data look like?

## 1

# Question Answering Research Areas

## Research Areas in Question Answering

Research Area	Details
<b>Knowledge-based QA (Semantic Parsing)</b>	<ul style="list-style-type: none"><li>• Answer is a logical form, possibly executed against a Knowledge Base</li><li>• Context is a Knowledge Base</li></ul>
Information Retrieval-based QA <ul style="list-style-type: none"><li>• Answer sentence selection</li><li>• Reading Comprehension</li></ul>	<ul style="list-style-type: none"><li>• Answer is a document, paragraph, sentence</li><li>• Context is a corpus of documents or a specific document</li></ul>
Visual QA	<ul style="list-style-type: none"><li>• Answer is simple and factual</li><li>• Context is one/multiple image(s)</li></ul>
Library Reference	<ul style="list-style-type: none"><li>• Answer is another question</li><li>• Context is the structured knowledge available in the library and the librarians' view of it.</li></ul>

## 2 Knowledge-based Question Answering

### Semantic Parsing

*Answering a natural language question by mapping it to a query over a structured database (formal representation of its meaning).*

*knowledge base*

*Question*

When was Justin Bieber born?

*Logical Form*

birth-year(Justin Bieber, x)

*KB Query*

*Knowledge base*

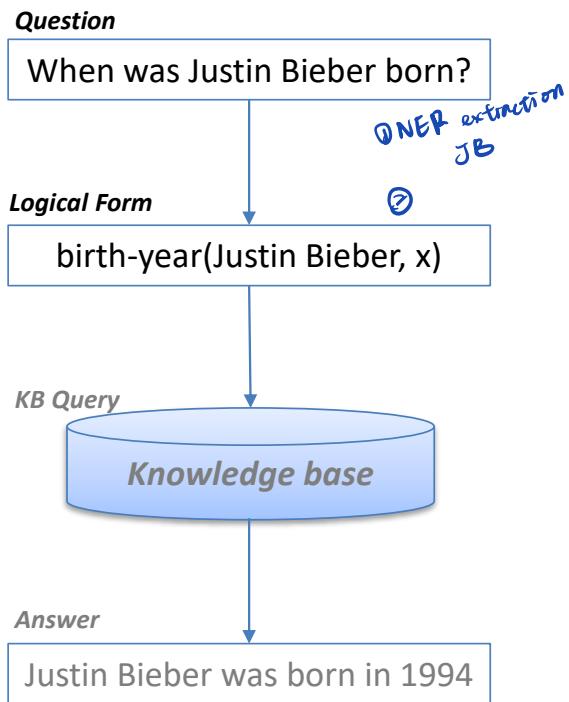
*Answer*

Justin Bieber was born in 1994

## 2 Knowledge-based Question Answering

### Semantic Parsing

*Answering a natural language question by mapping it to a query over a structured database (formal representation of its meaning).*



*Mapping from a text string to any logical form*

Question	Logical Form
When was Justin Bieber born?	birth-year(Justin Bieber, x)
What is the largest state?	$\text{argmax}(\lambda x.\text{state}(x), \lambda x.\text{size}(x))$

**How to map?** Map either to some version of predicate calculus or a query language like SQL or SPARQL

<https://query.wikidata.org/>

**Use expensive supervised data?**

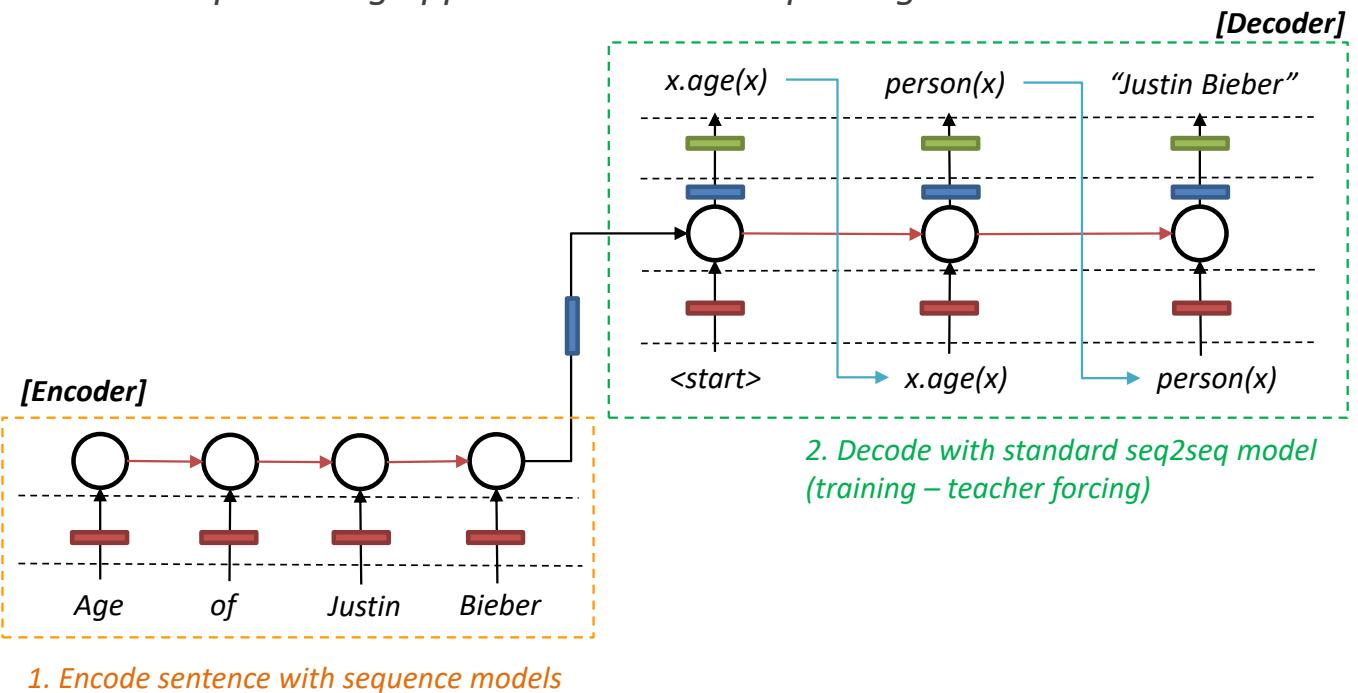
Require experts for the manual annotation process....

## 2 Knowledge-based Question Answering

### Seq2Seq model for semantic parser

*How to transfer the text to the logical form?*

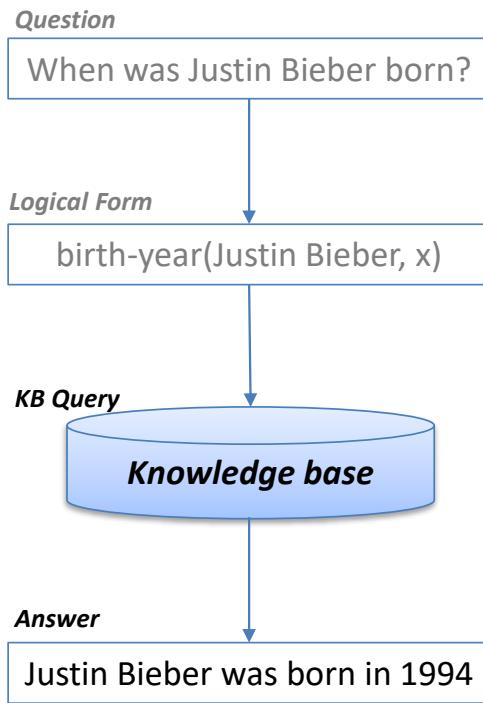
*A basic deep learning approach to semantic parsing*



## 2 Knowledge-based Question Answering

### Semantic Parsing

*Answering a natural language question by mapping it to a query over a structured database (formal representation of its meaning).*



- ① search from knowledge base  
 ② generate answer based on human language

*Answer questions that ask about one of the missing arguments in a triple*

Subject	Predicate (relation)	Object
Justin Bieber	birth-year	1994
Frédéric Chopin	birth-year	1810
...	...	...

wikipedia

- DBpedia
- Freebase

**How to produce the answer?**

- Seq2seq
- Template based generation

## 2 Knowledge-based Question Answering

### Pros and Cons of Knowledge-based QA

- *Logical Form instead of (direct) answer makes system robust*
- *Answer independent of question and parsing mechanism*
- *Constrained to queriable questions in Database Schema*
- *Difficult to find the well-structured training dataset*

# 0 Question Answering

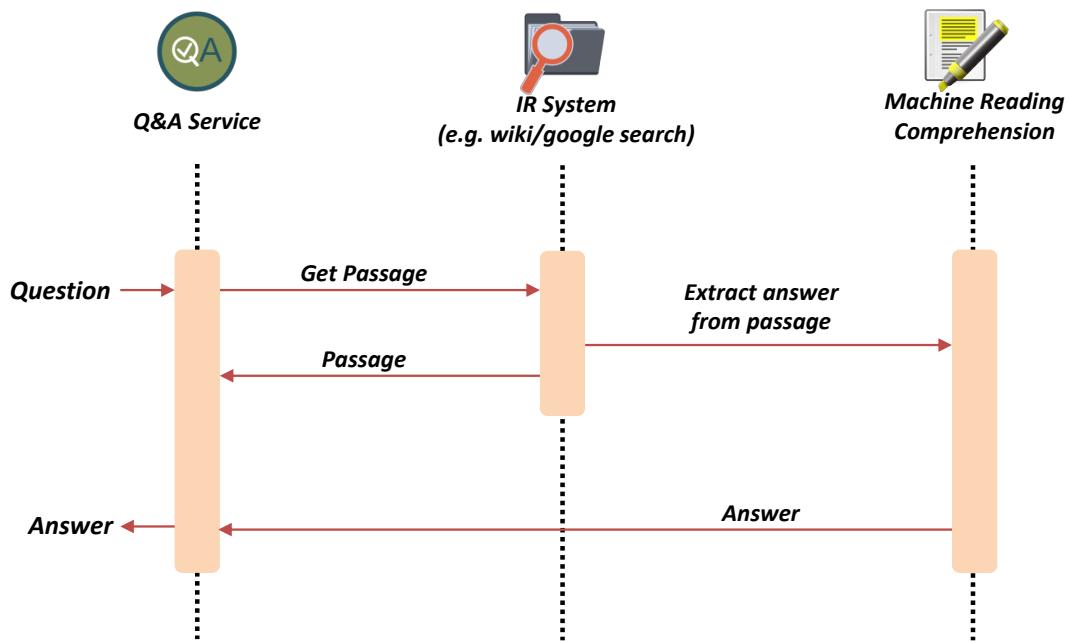
## Research Areas in Question Answering

Research Area	Details
Knowledge-based QA (Semantic Parsing)	<ul style="list-style-type: none"><li>• Answer is a logical form, possibly executed against a Knowledge Base</li><li>• Context is a Knowledge Base</li></ul>
<b>Information Retrieval-based QA</b> <ul style="list-style-type: none"><li>• Answer sentence selection</li><li>• Reading Comprehension</li></ul>	<ul style="list-style-type: none"><li>• Answer is a document, paragraph, sentence</li><li>• Context is a corpus of documents or a specific document</li></ul>
Visual QA	<ul style="list-style-type: none"><li>• Answer is simple and factual</li><li>• Context is one/multiple image(s)</li></ul>
Library Reference	<ul style="list-style-type: none"><li>• Answer is another question</li><li>• Context is the structured knowledge available in the library and the librarians' view of it.</li></ul>

### 3 Information Retrieval-based Question Answering

#### Information Retrieval-based Question Answering

Answering a user's question by *finding short text segments, sentences, or documents* on the web or collection of document



# 3 Reading Comprehension

## Information Retrieval-based Question Answering

*Answering a user's question by finding short text segments, sentences, or documents on the web or collection of document*

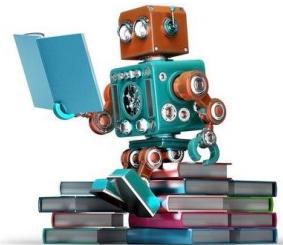
- *Reading Comprehension and Answer Sentence Selection:*
  - *Finding an answer in a paragraph or a document*
  - *Picking a suitable sentence from a corpus that can be used to answer a question*

# 3 Reading Comprehension

## Reading Comprehension

*To answer these questions, you need to first gather information by collecting answer-related sentences from the article.*

***Can we teach this to machine?***



***Yes, we can!***

*Machine Comprehension of Text*  
(Burges 2013)

### THE BOAT PARADE

The boats are floating along the lakeshore. It is the summer boat parade.

There are motor boats, rowboats and sailboats.

Jessica's favorite is the yellow motor boat with the flag. The rowboat decorated with flowers is Lisa's favorite. Tony likes the purple sailboat.

The boats float by one at a time. The people on the boats waive at the crowds. The crowds cheer the boats.

The boat parade is so much fun to watch. It is the best part of the summer.

**Answer the Questions:**

1. Where are the boats floating?
2. What kind of boats are there?
3. What is Lisa's favorite boat?

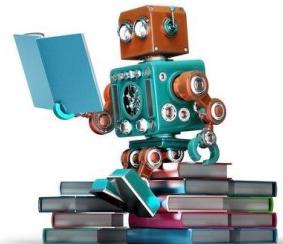


### 3 Reading Comprehension

## Reading Comprehension

*To answer these questions, you need to first gather information by collecting answer-related sentences from the article.*

***Can we teach this to machine?***



*A machine comprehends a passage of text if, for any question regarding that text that can be answered correctly by a majority of native speakers*

#### THE BOAT PARADE

The boats are floating along the lakeshore. It is the summer boat parade.

There are motor boats, rowboats and sailboats.

Jessica's favorite is the yellow motor boat with the flag. The rowboat decorated with flowers is Lisa's favorite. Tony likes the purple sailboat.

The boats float by one at a time. The people on the boats waive at the crowds. The crowds cheer the boats.

The boat parade is so much fun to watch. It is the best part of the summer.

**Answer the Questions:**

1. Where are the boats floating?
2. What kind of boats are there?
3. What is Lisa's favorite boat?

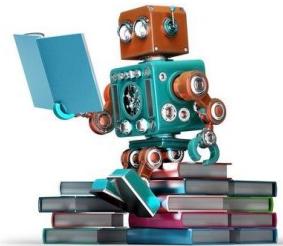


### 3 Reading Comprehension

## Reading Comprehension

*To answer these questions, you need to first gather information by collecting answer-related sentences from the article.*

### Why do we need to teach this?



**The ability to comprehend text will lead us to a better search and solve lots of NLP problems!**

#### THE BOAT PARADE

The boats are floating along the lakeshore. It is the summer boat parade.

There are motor boats, rowboats and sailboats.

Jessica's favorite is the yellow motor boat with the flag. The rowboat decorated with flowers is Lisa's favorite. Tony likes the purple sailboat.

The boats float by one at a time. The people on the boats waive at the crowds. The crowds cheer the boats.

The boat parade is so much fun to watch. It is the best part of the summer.

#### Answer the Questions:

1. Where are the boats floating?
2. What kind of boats are there?
3. What is Lisa's favorite boat?



# 3 Reading Comprehension

## Corpora for Reading Comprehension

Dataset	Answer Type	Domain
MCTest (Richardson et al. 2013)	Multiple choice	Children's stories
CNN/Daily Mail (Hermann et al. 2015)	Spans	News
Children's book test (Hill et al. 2016)	Multiple choice	Children's stories
SQuAD (Rajpurkar et al., 2016)	Spans	Wikipedia
MS MARCO (Nguyen et al., 2016)	Free-from text, Unanswerable	Web Search
NewsQA (Trischler et al., 2017)	Spans	News
SearchQA (Dunn et al., 2017)	Spans	Jeopardy
TriviaQA (Joshi et al., 2017)	Spans	Trivia
RACE (Lai et al., 2017)	Multiple choice	Mid/High School Exams
Narrative QA (Kočiský et al., 2018)	Free-form text	Movie Scripts, Literature
SQuAD 2.0 (Rajpurkar et al., 2018)	Spans, Unanswerable	Wikipedia

### 3 Reading Comprehension

## TriviaQA: A Large Scale Dataset for Reading Comprehension

### TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension

The full dataset is coming soon. Here's a sneak peek! The evidence documents come from two domains -- Wikipedia and the web. Click on the "Evidence" button to see the document for each question.

QuestionId	Question	Answer	Web	Wikipedia
qw_3199	Miami Beach in Florida borders which ocean?	Atlantic	<a href="#">Evidence</a>	<a href="#">Evidence</a>
bt_1255	What was the occupation of Lovely Rita according to the song by the Beatles	Traffic Warden	<a href="#">Evidence</a>	<a href="#">Evidence</a>
qg_77	Who was Poopdeck Pappys most famous son?	Popeye	<a href="#">Evidence</a>	<a href="#">Evidence</a>
wh_1026	The Nazi regime was Germany's Third Reich; which was the first Reich?	HOLY ROMAN EMPIRE	<a href="#">Evidence</a>	<a href="#">Evidence</a>
bb_1342	At which English racecourse did two horses collapse and die in the parade ring due to electrocution, in February 2011?	Newbury	<a href="#">Evidence</a>	<a href="#">Evidence</a>
wh_2759	Which type of hat takes its name from an 1894 novel by George Du Maurier where the title character has the surname O'Ferrall ?	TRILBY	<a href="#">Evidence</a>	<a href="#">Evidence</a>
sfq_8522	What was the Elephant Man's real name?	Joseph Merrick	<a href="#">Evidence</a>	<a href="#">Evidence</a>

### 3 Reading Comprehension

## TriviaQA: A Large Scale Dataset for Reading Comprehension

Our UsydNLP achieved the **No.1** in the TriviaQA Leaderboard (Web Setting)!

[Wikipedia](#)
[Web](#)

Phase description

This phase refers to the Web domain described in the paper.

Max submissions per day: 3

Max submissions total: 100

[Download CSV](#)

Results								
#	User	Entries	Date of Last Entry	Team Name	full-em ▲	full-ft ▲	verified-em ▲	verified-ft ▲
1	<a href="#">mandarjoshi</a>	10	12/02/19	<b>Dataset Author (Oracle Result)</b>	82.99 (1)	87.18 (1)	90.38 (1)	92.96 (1)
2	<a href="#">usydnlp</a>	11	01/11/21		72.17 (2)	77.36 (2)	84.40 (2)	87.11 (2)
3	<a href="#">NEUKG</a>	9	08/12/19		69.64 (3)	73.80 (3)	83.36 (3)	85.66 (4)
4	<a href="#">mingyan</a>	1	02/26/18	SLQA	68.65 (4)	73.07 (5)	82.44 (5)	85.35 (5)
5	<a href="#">scv.back</a>	4	06/07/18	S3R	68.21 (5)	73.26 (4)	82.57 (4)	86.05 (3)
6	<a href="#">dirkweissenborn</a>	7	01/15/18		67.46 (6)	72.80 (6)	77.63 (9)	82.01 (8)
7	<a href="#">S3R</a>	5	02/28/18		66.82 (7)	71.91 (7)	81.01 (6)	84.12 (6)
8	<a href="#">chrisc</a>	3	09/24/17		66.37 (8)	71.32 (8)	79.97 (7)	83.70 (7)
9	<a href="#">Mary</a>	5	05/26/20		66.09 (9)	69.78 (9)	79.71 (8)	81.88 (9)
10	<a href="#">shuohang</a>	3	11/13/17		63.04 (10)	68.53 (10)	69.70 (10)	74.57 (10)
11	<a href="#">swabha</a>	1	10/27/17	swabha_ankur_tom	53.75 (11)	58.57 (11)	63.20 (11)	66.88 (11)
12	<a href="#">hux444</a>	2	07/25/17		46.65 (12)	52.89 (12)	56.96 (12)	61.48 (12)

<http://nlp.cs.washington.edu/triviaqa/sample.html>

### 3 Reading Comprehension

## SQuAD: Stanford Question Answering Dataset

### Victoria\_(Australia)

#### The Stanford Question Answering Dataset

The economy of Victoria is highly diversified: service sectors including financial and property services, health, education, wholesale, retail, hospitality and manufacturing constitute the majority of employment. Victoria's total gross state product (GSP) is ranked second in Australia, although Victoria is ranked fourth in terms of GSP per capita because of its limited mining activity. Culturally, Melbourne is home to a number of museums, art galleries and theatres and is also described as the "sporting capital of Australia". The Melbourne Cricket Ground is the largest stadium in Australia, and the host of the 1956 Summer Olympics and the 2006 Commonwealth Games. The ground is also considered the "spiritual home" of Australian cricket and Australian rules football, and hosts the grand final of the Australian Football League (AFL) each year, usually drawing crowds of over 95,000 people. Victoria includes eight public universities, with the oldest, the University of Melbourne, having been founded in 1853.

**What kind of economy does Victoria have?**

*Ground Truth Answers:* diversified highly

diversified highly diversified

*Prediction:* highly diversified

**Where according to gross state product does Victoria rank in Australia?**

*Ground Truth Answers:* second second second

*Prediction:* second

**At what rank does GPS per capita set Victoria?**

*Ground Truth Answers:* fourth fourth fourth

*Prediction:* fourth

**What city in Victoria is called the sporting capital of Australia?**

*Ground Truth Answers:*

Melbourne Melbourne Melbourne

*Prediction:* Melbourne

### 3 Reading Comprehension

#### A Generic Neural Model for Reading Comprehension

*Step 1: For both **documents** and **questions**, convert words to word vectors*

*span-based*



##### Document (*D*)

A partly submerged glacier cave on Perito Moreno Glacier. The ice facade is approximately 60 m high. Ice formations in the Titlis glacier cave. A glacier cave is a cave formed within the ice of a glacier. Glacier caves are often called ice caves, but the latter term is properly used to describe bedrock caves that contain year-round ice

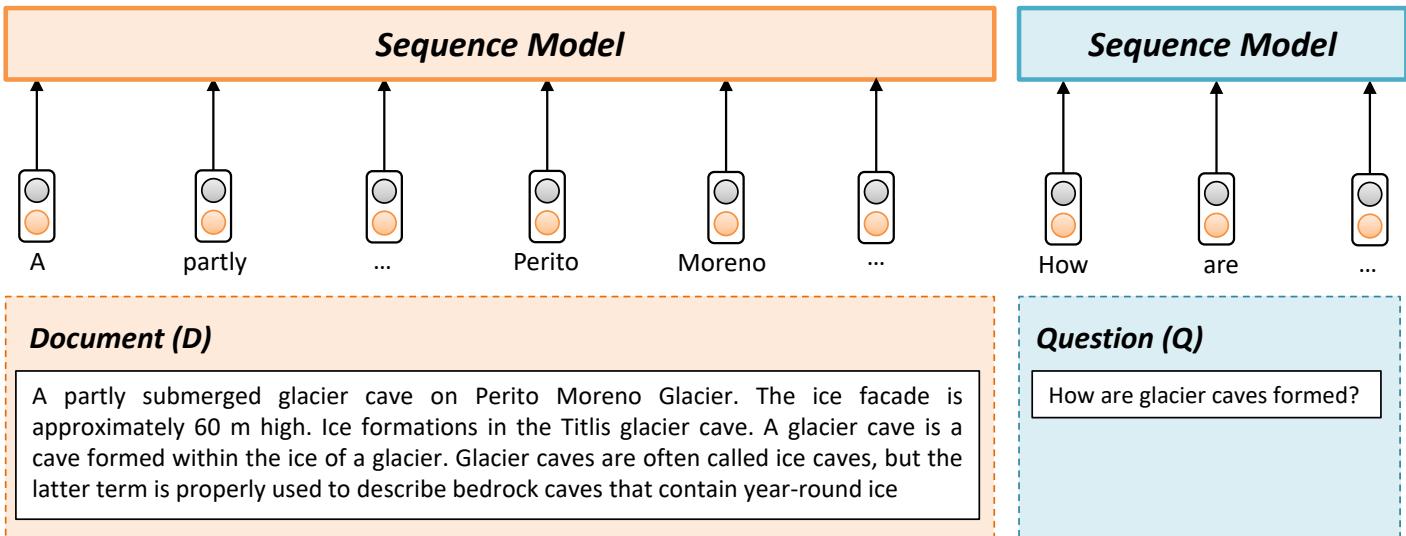
##### Question (*Q*)

How are glacier caves formed?

### 3 Reading Comprehension

## A Generic Neural Model for Reading Comprehension

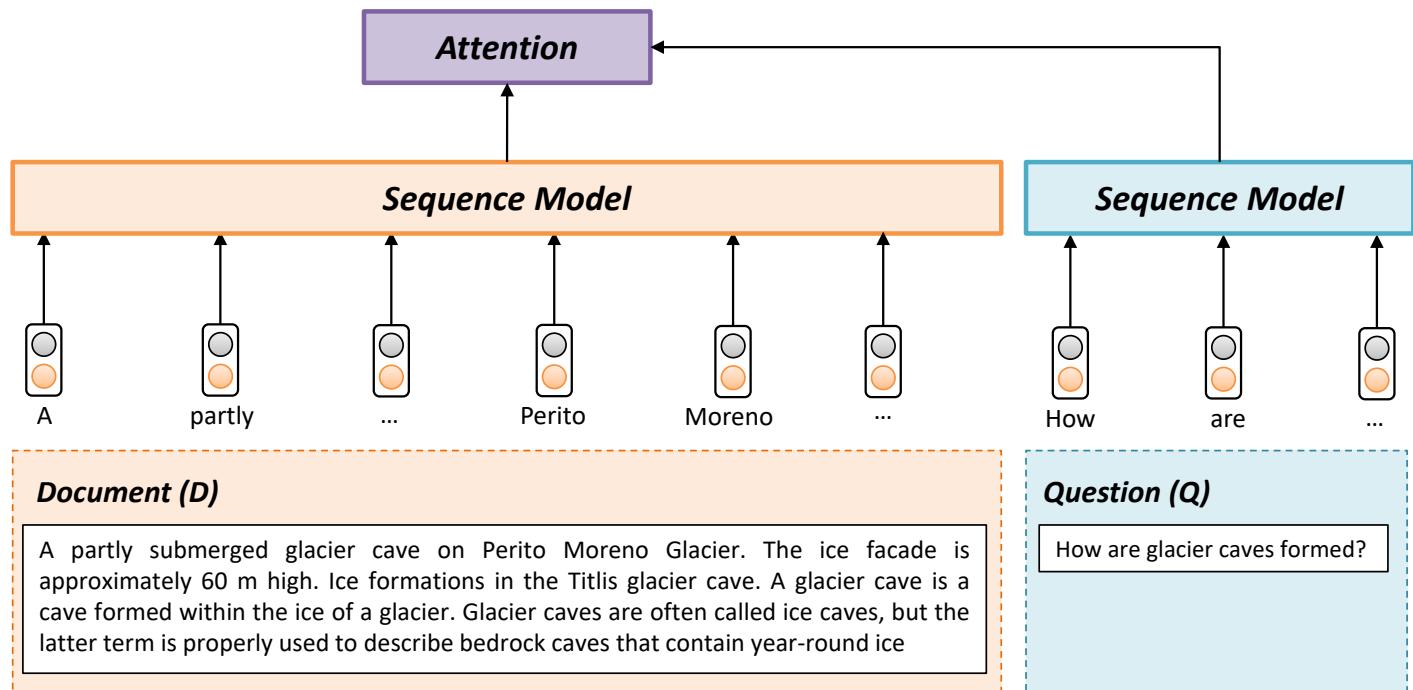
*Step2: Encode context (documents) and question with sequence models*



### 3 Reading Comprehension

#### A Generic Neural Model for Reading Comprehension

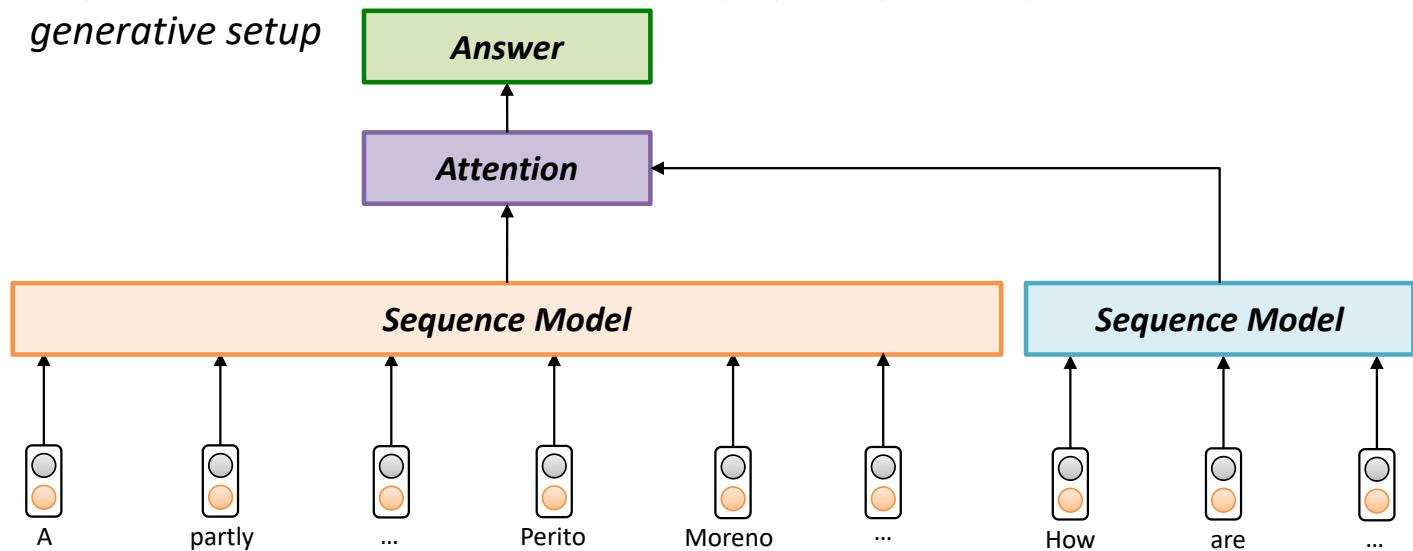
*Step3: Combine context (documents) and question with an attention*



### 3 Reading Comprehension

#### A Generic Neural Model for Reading Comprehension

**Step4:** Select *answer* from attention map by using a classifier or with generative setup



#### Document (D)

A partly submerged glacier cave on Perito Moreno Glacier. The ice facade is approximately 60 m high. Ice formations in the Titlis glacier cave. A glacier cave is a cave formed within the ice of a glacier. Glacier caves are often called ice caves, but the latter term is properly used to describe bedrock caves that contain year-round ice

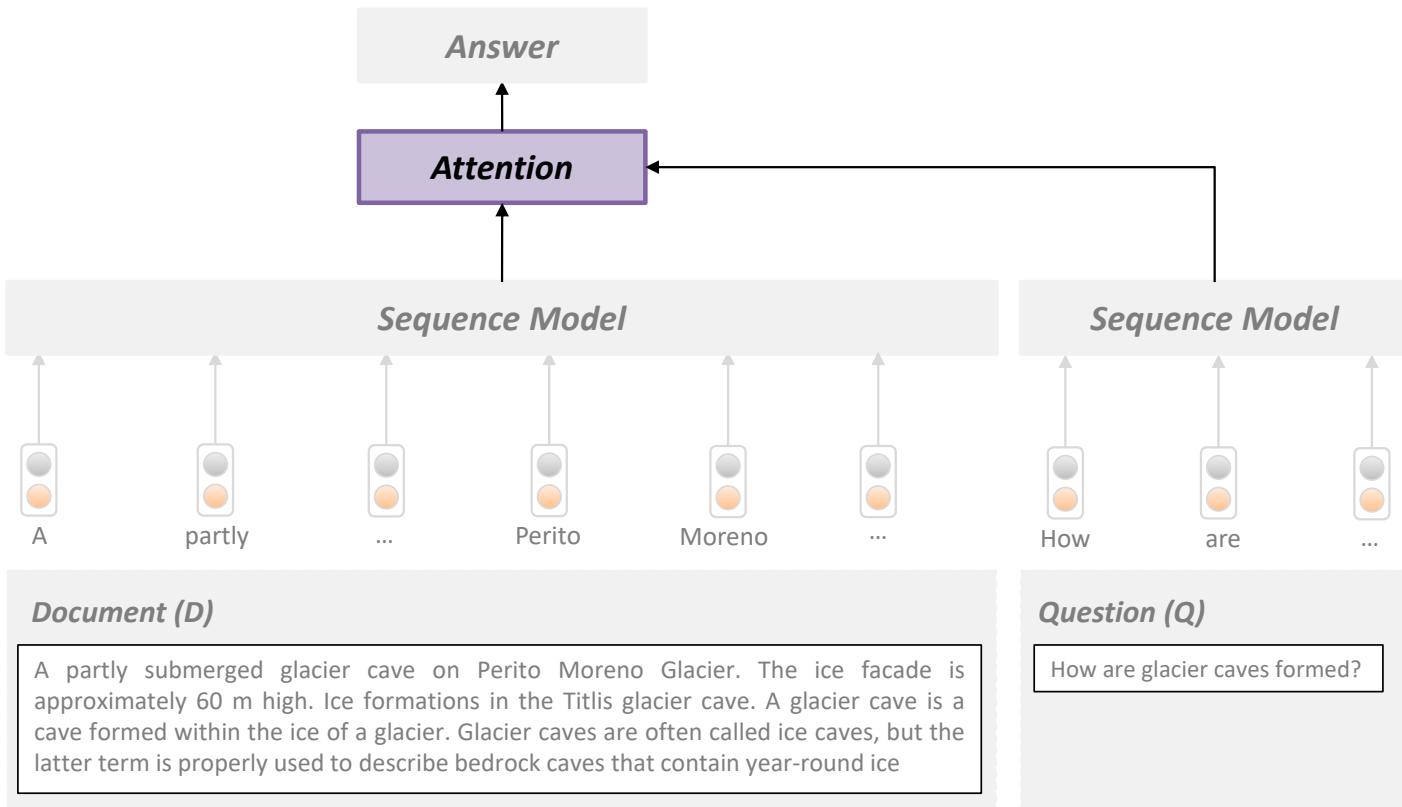
#### Question (Q)

How are glacier caves formed?

### 3 Reading Comprehension

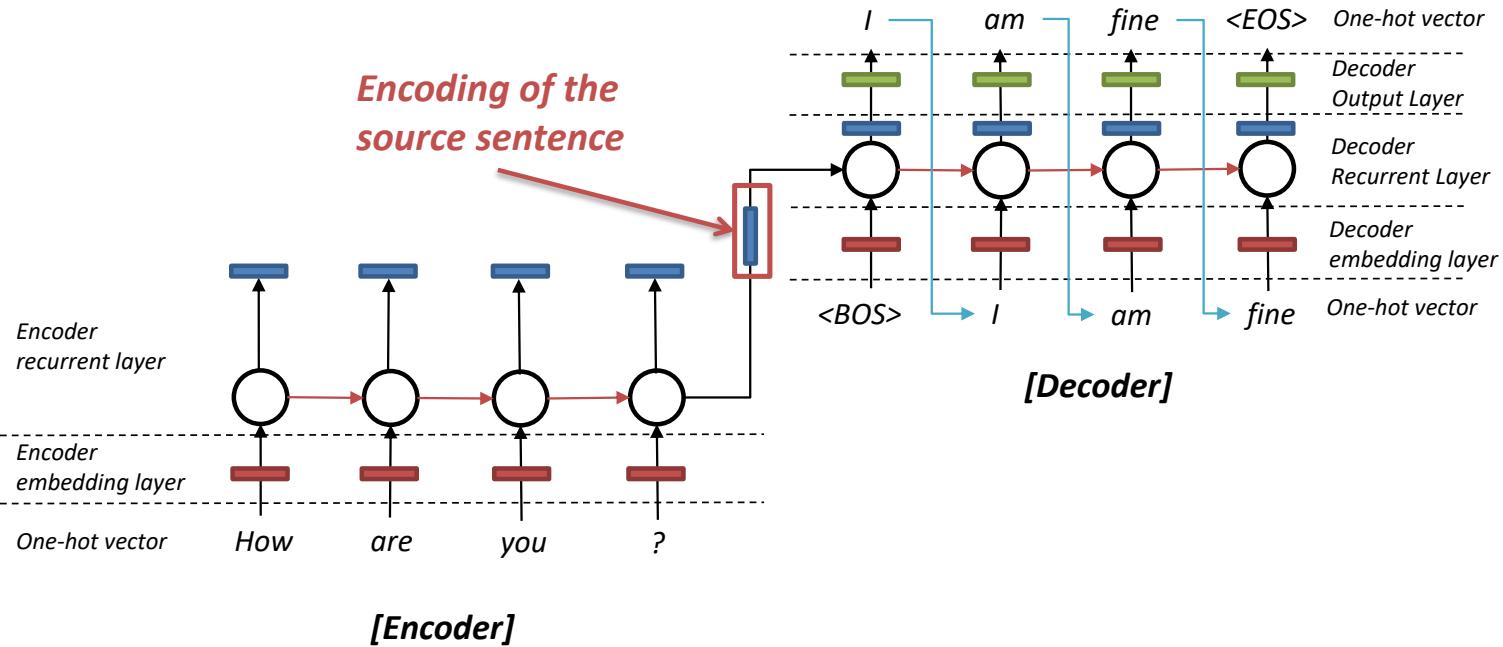
#### A Generic Neural Model for Reading Comprehension

*What is the Attention? Why we need this?*



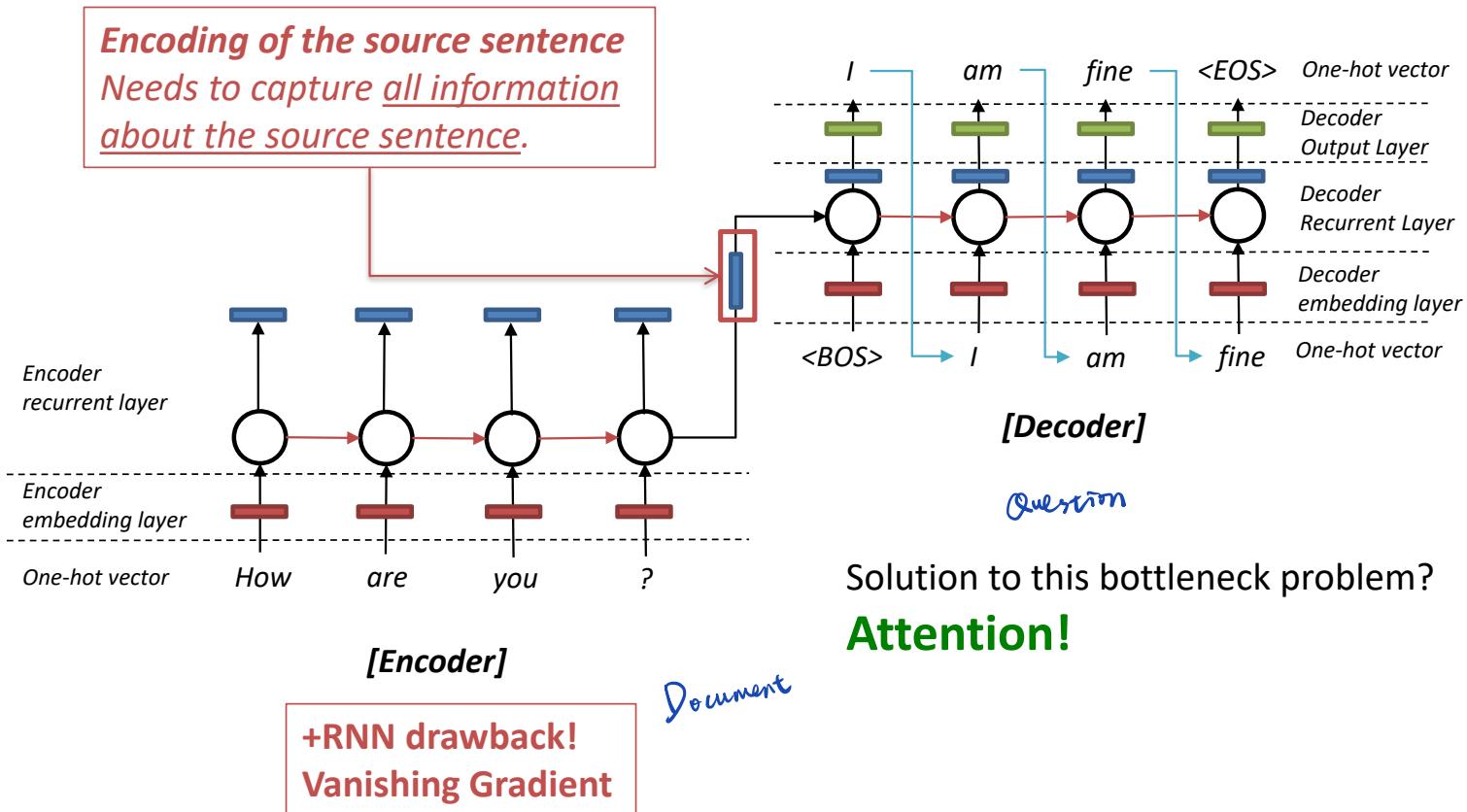
# 4 Attention

## Seq2Seq Model: Recap



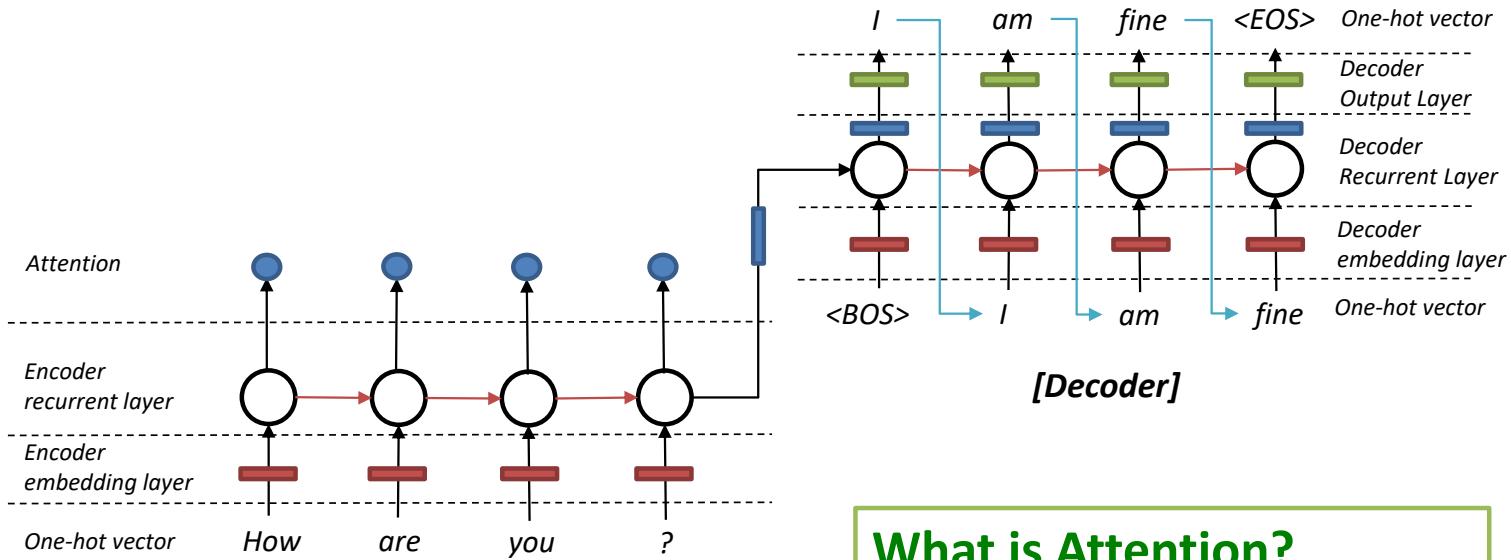
## 4 Attention

### Seq2Seq Model: the bottleneck problem



# 4 Attention

## Seq2Seq with Attention

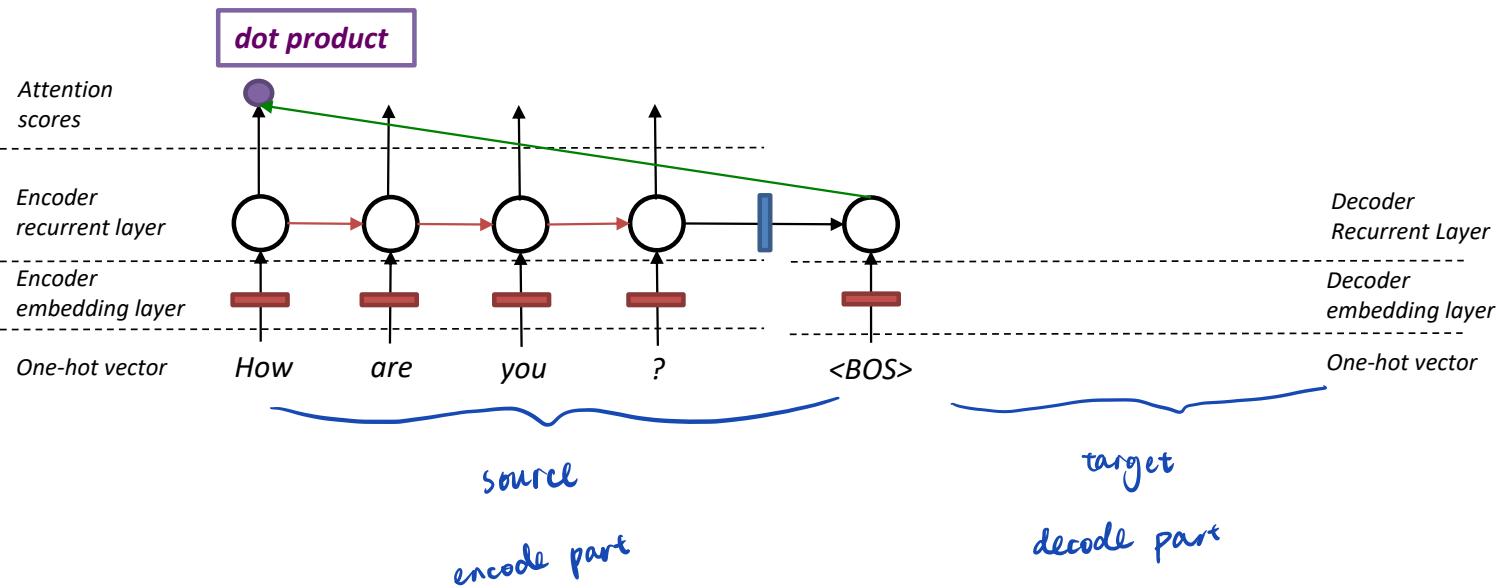


**What is Attention?**

*On each step of the decoder, use direct connection to the encoder to focus on a particular part of the input sequence*

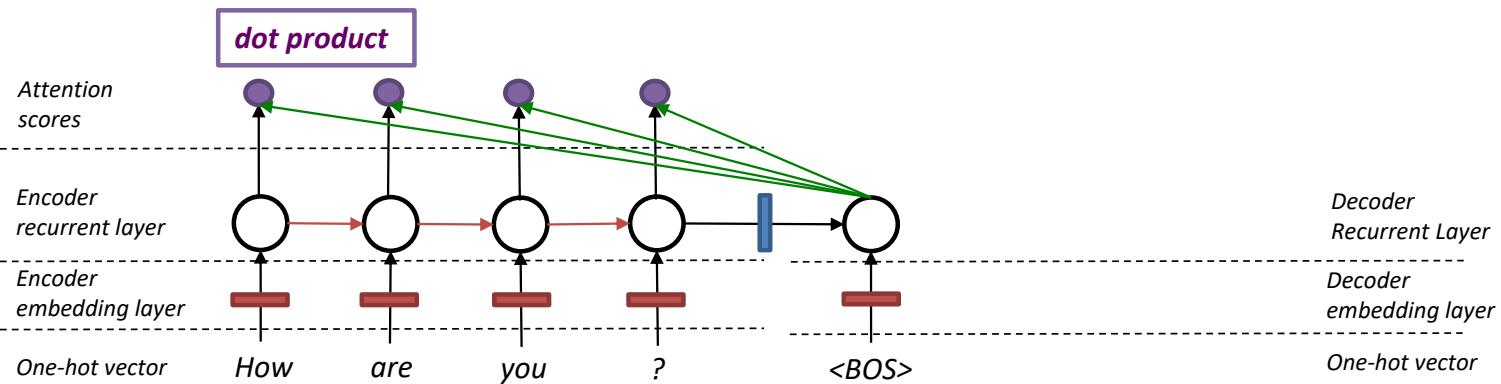
## 4 Attention

### Seq2Seq with Attention



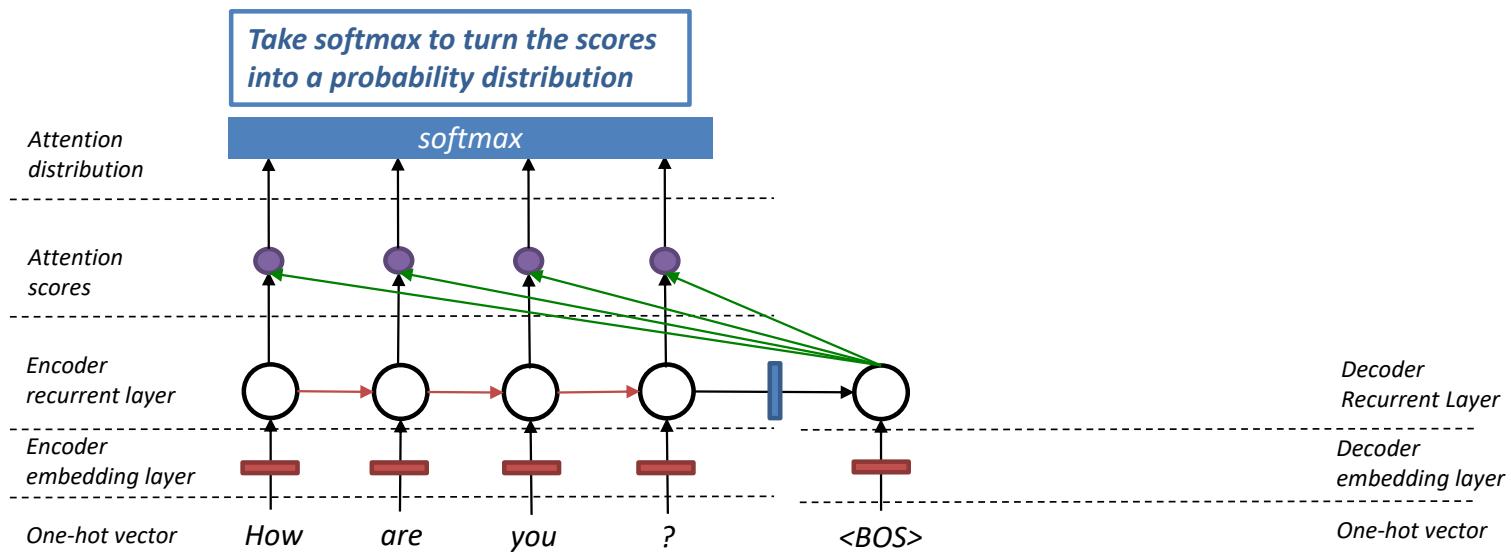
## 4 Attention

### Seq2Seq with Attention



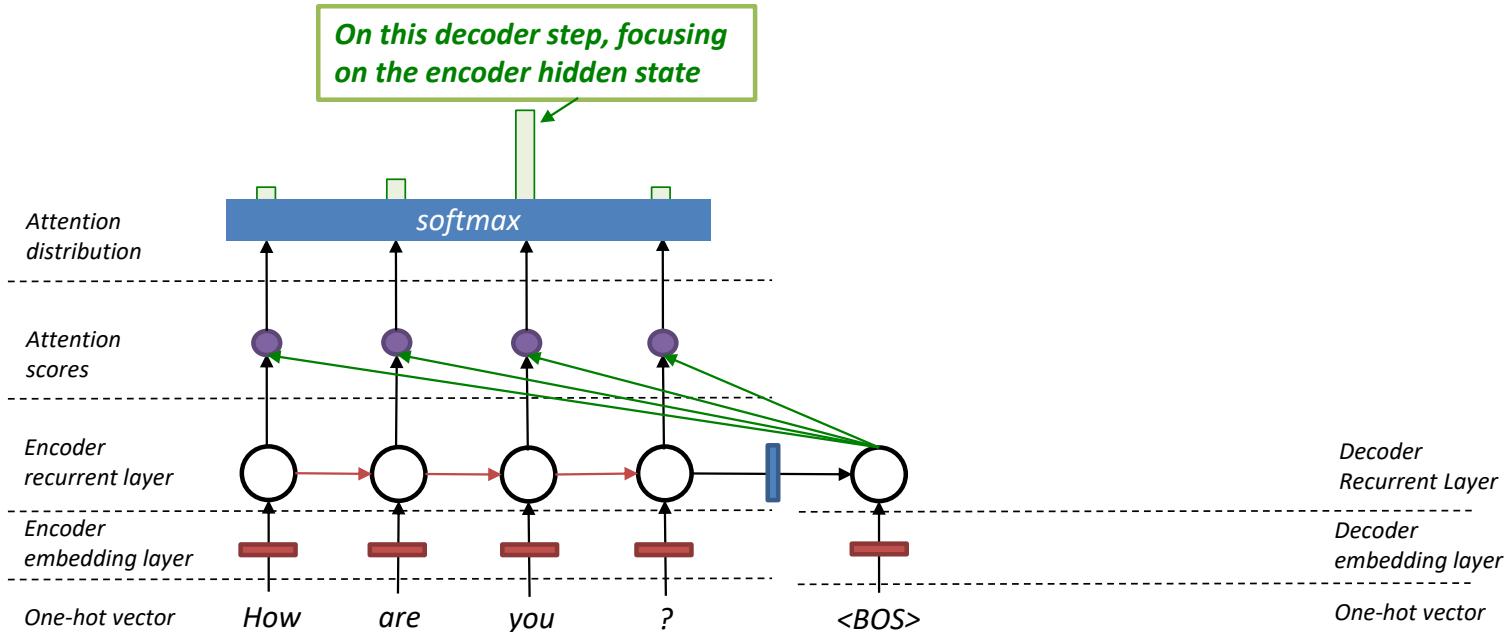
## 4 Attention

### Seq2Seq with Attention



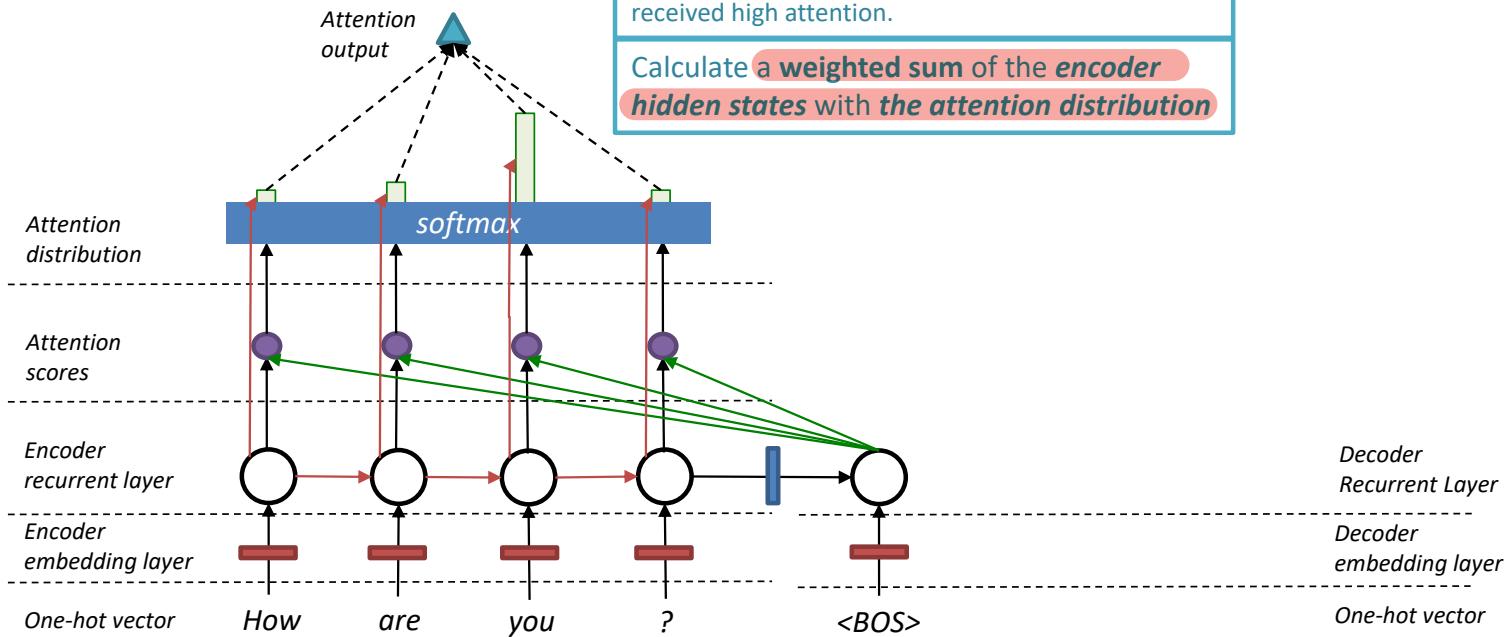
# 4 | Attention

# Seq2Seq with Attention

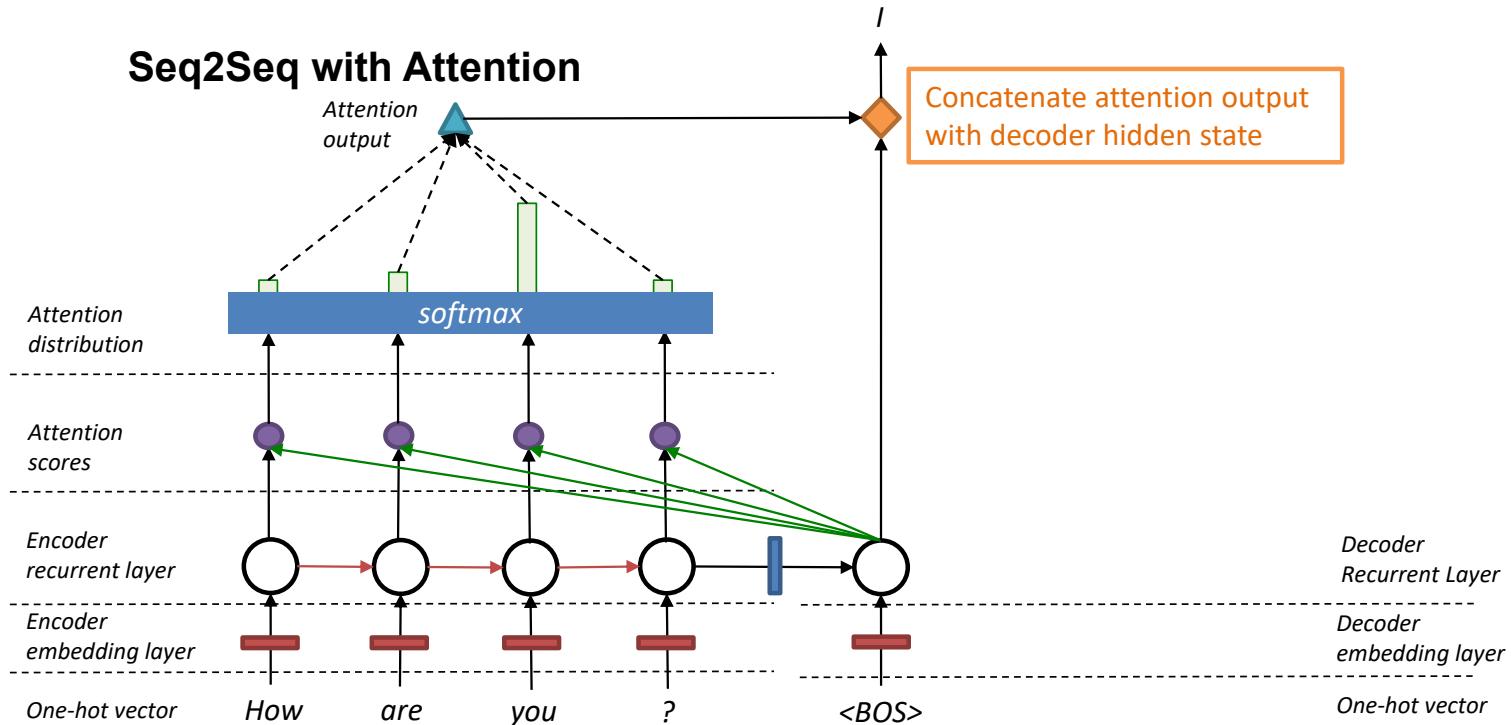


# 4 Attention

## Seq2Seq with Attention

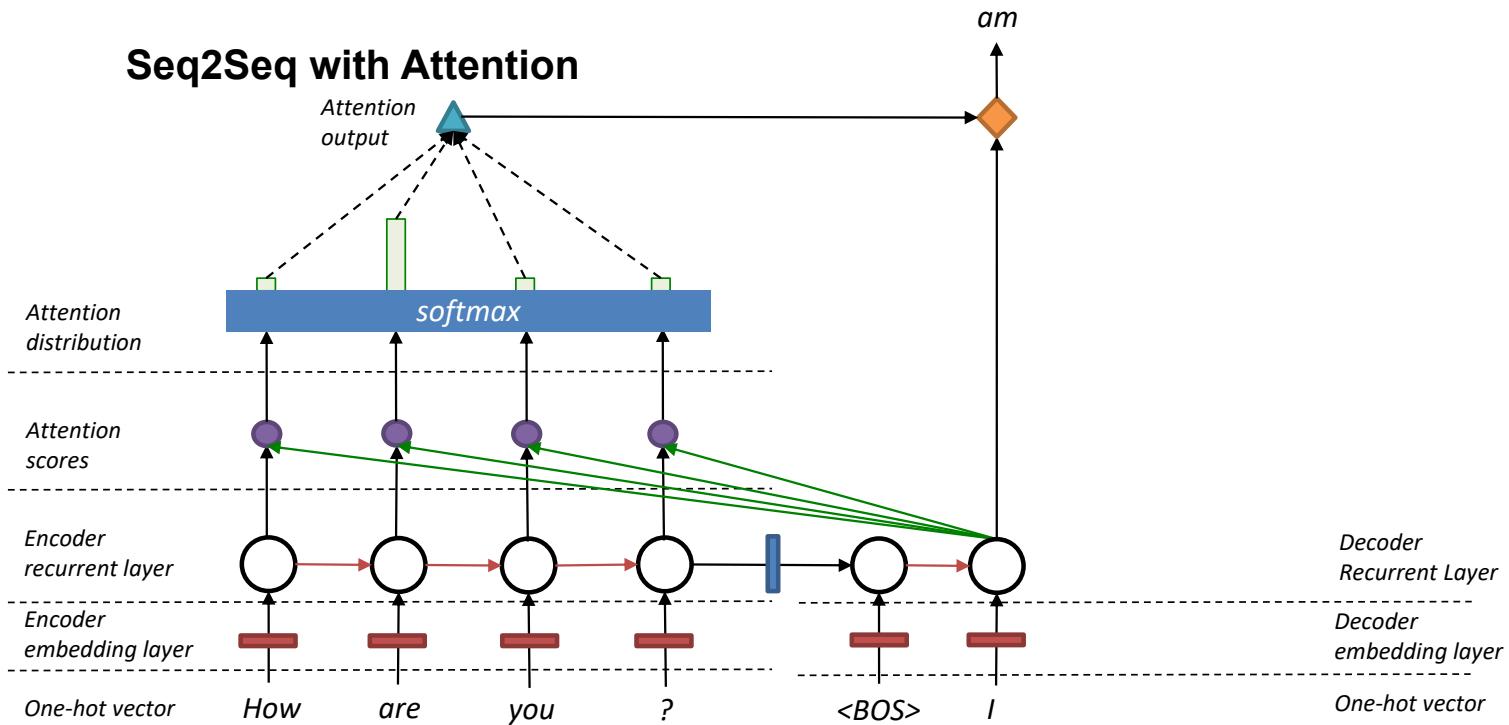


## 4 Attention



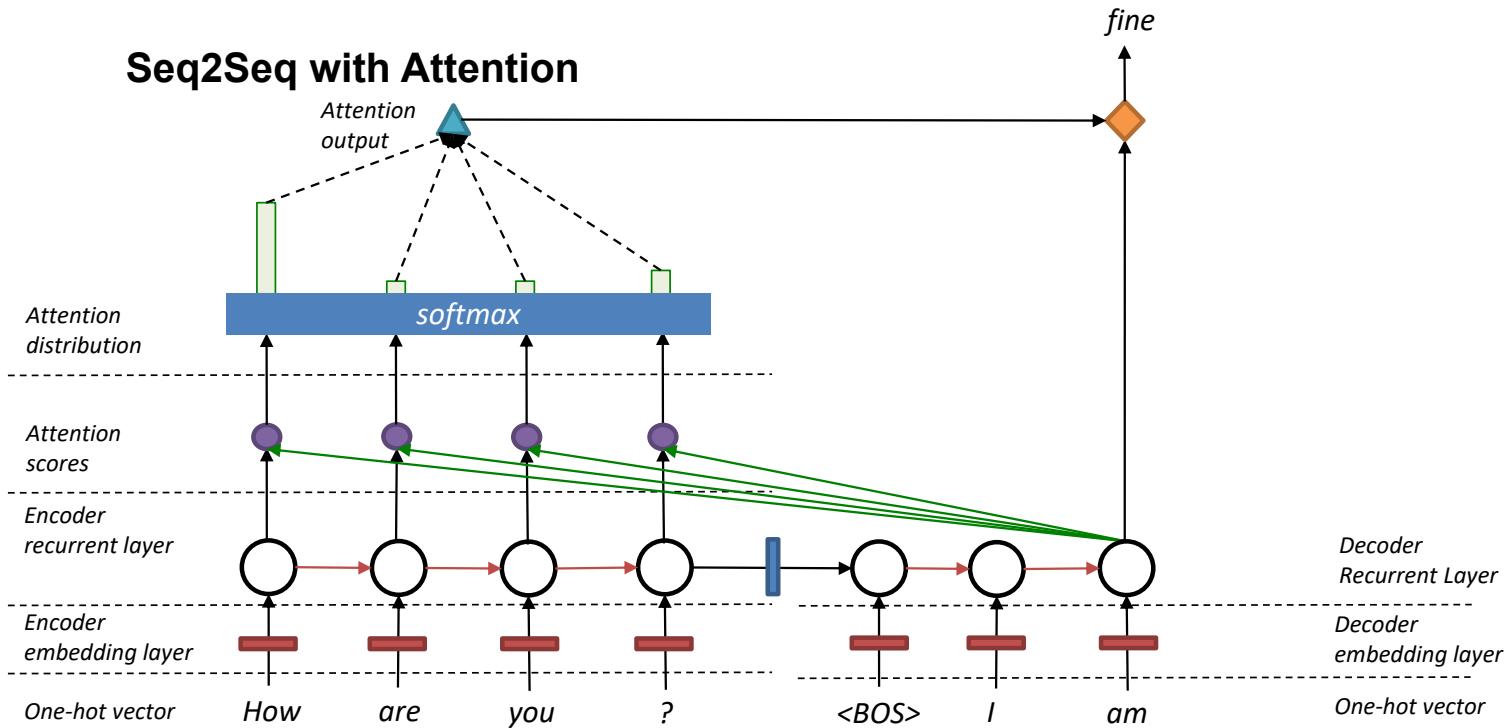
## 4 Attention

### Seq2Seq with Attention



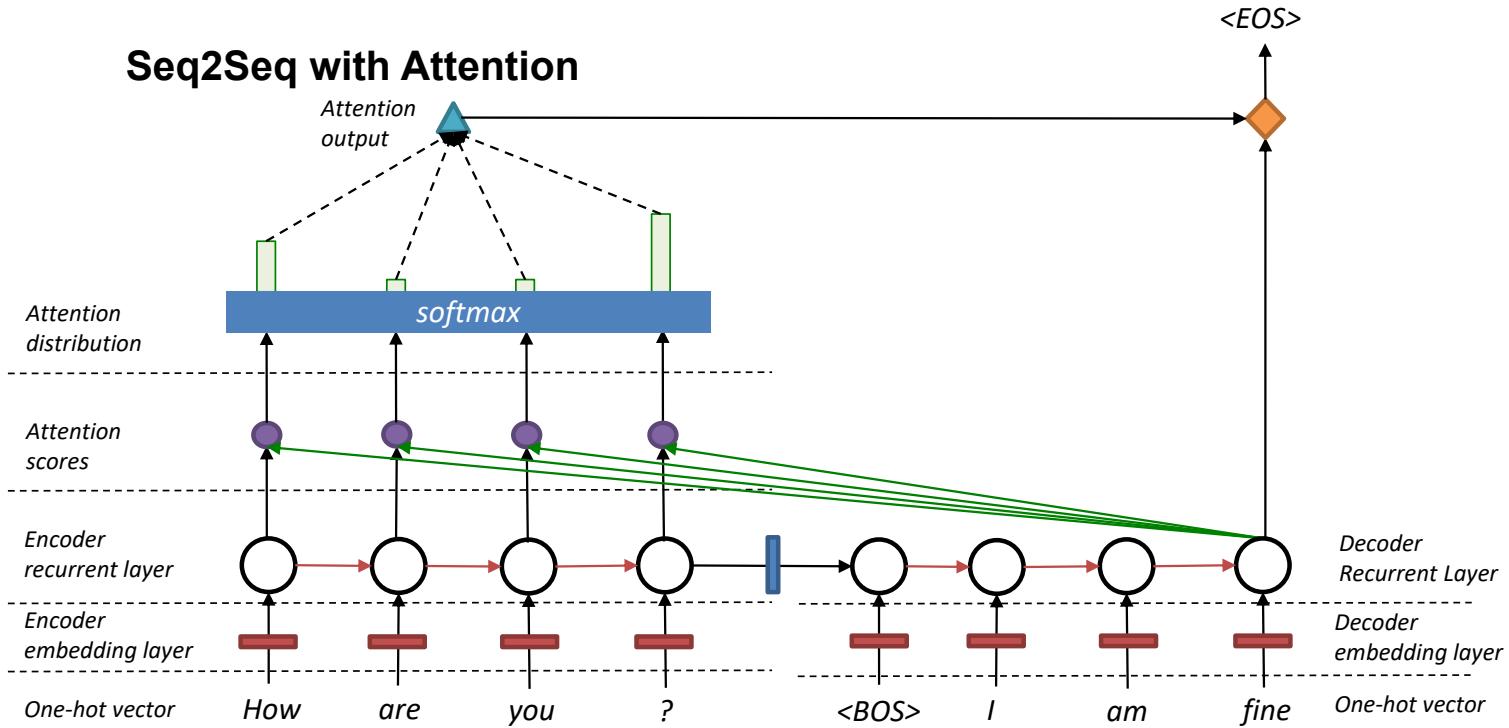
## 4 Attention

### Seq2Seq with Attention



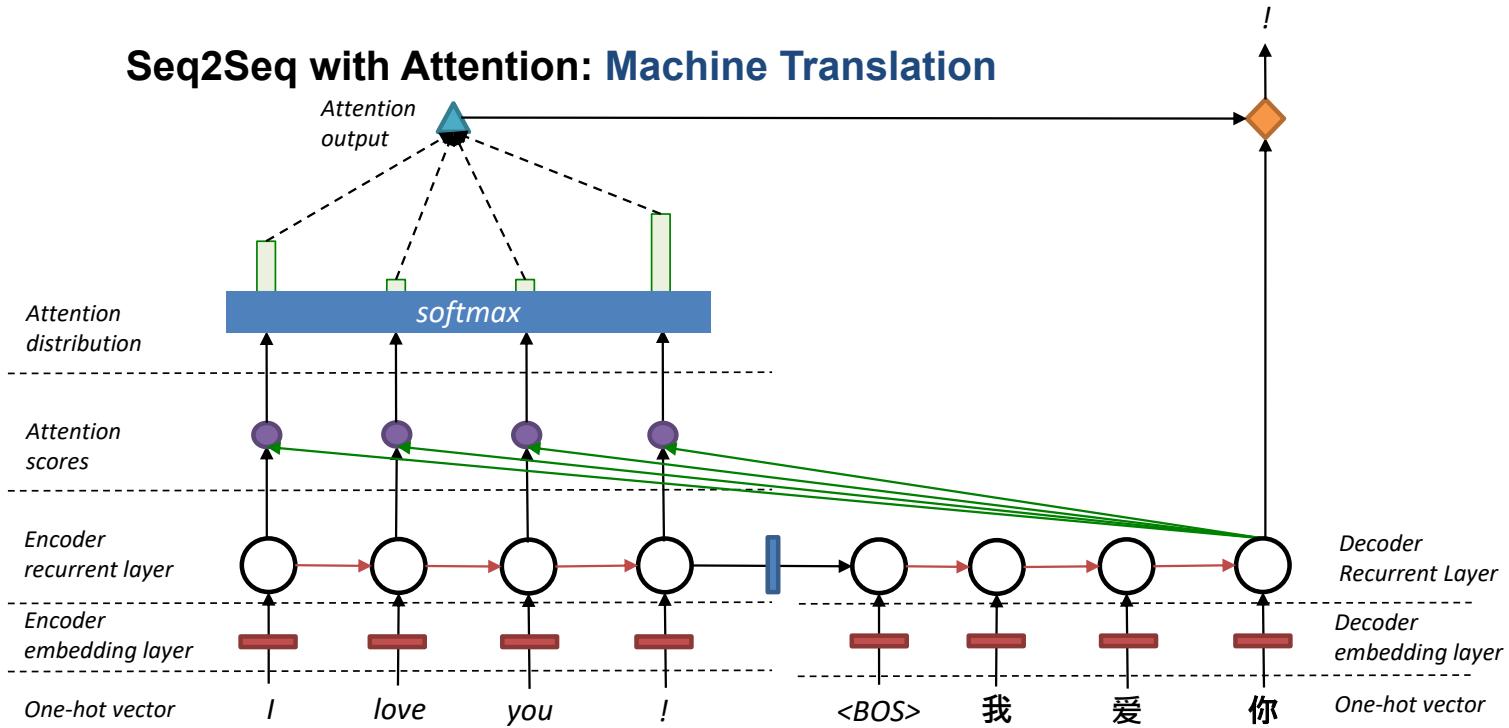
# 4 Attention

## Seq2Seq with Attention



## 4 Attention

### Seq2Seq with Attention: Machine Translation



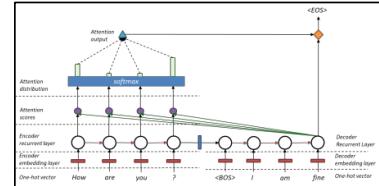
## 4 Attention

### Seq2Seq with Attention (Equations)

- *Encoder hidden states:*  $h_1, \dots, h_N \in \mathbb{R}^h$
  - *Decoder hidden state:*  $s_t \in \mathbb{R}^h$  (on timestep  $t$ )
1. *Attention score*  $\mathbf{e}^t : \mathbf{e}^t = [\mathbf{s}_t^T \mathbf{h}_1, \dots, \mathbf{s}_t^T \mathbf{h}_N] \in \mathbb{R}^N$  (for timestep  $t$ )
  2. *Use softmax to get the attention distribution,  $\alpha^t$  (for timestep  $t$ )*  
(this is a probability distribution and sums to 1)
 
$$\alpha^t = \text{softmax}(\mathbf{e}^t) \in \mathbb{R}^N$$
  3. *Attention Output: Use  $\alpha^t$  to take a weighted sum of the encoder hidden states*

$$\mathbf{a}_t = \sum_{i=1}^N \alpha_i^t \mathbf{h}_i \in \mathbb{R}^h$$
  4. *Then, concatenate the attention output  $\mathbf{a}_t$  with the decoder hidden state  $s_t$  and proceed as in the non-attention seq2seq model*

$$[\mathbf{a}_t; \mathbf{s}_t] \in \mathbb{R}^{2h}$$



## 4 Attention

### Why we use Attention? The benefit!

#### *Improve performance*

- Allow decoder to focus on certain parts of the source

#### *Solving the bottleneck problem*

- Allow decoder to directly look at the source (input)

#### *Reducing vanishing gradient problem*

- Provide shortcut to faraway states

#### *Providing some interpretability*

- Inspect attention distribution, and show what the decoder was focusing on

## 4 Attention

## Attention is now a general component in Deep Learning NLP

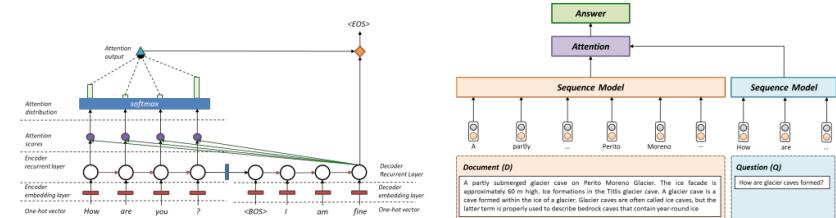
*Attention is great way to improve the sequence to sequence model.*

You can use attention in many architecture (not just seq2seq) and many NLP tasks (not just dialog system/NLG, Translation)

### ***More general definition of attention:***

*Given a set of vector values, and a vector query, attention is a technique to compute a weighted sum of the values, dependent on the query.*

For example, in the seq2seq + attention model, each decoder hidden state (**query**) attends to all the encoder hidden states (**values**).



## 4 Attention

### Attention variants

*There are several ways to compute attention score.*

- Encoder hidden states:  $h_1, \dots, h_N \in \mathbb{R}^h$
- Decoder hidden state:  $s_t \in \mathbb{R}^h$  (on timestep  $t$ )

Attention Name	Attention score function	Reference
<b>Content-base</b>	$score(s_t, h_i) = \text{cosine}[s_t, h_i]$	<a href="#">Graves 2014</a>
<b>Dot-product</b>	$score(s_t, h_i) = s_t^\top h_i$	<a href="#">Luong 2015</a>
<b>Scaled Dot-product</b>	$score(s_t, h_i) = \frac{s_t^\top h_i}{\sqrt{n}}$ *NOTE: very similar to the dot-product attention except for a scaling factor; where $n$ is the dimension of the source hidden state.	<a href="#">Vaswani 2017</a>
<b>Additive</b>	$score(s_t, h_i) = v_a^\top \tanh(W_a[s_t; h_i])$	<a href="#">Vaswani 2017</a>
<b>General</b>	$score(s_t, h_i) = s_t^\top W_d h_i$ *NOTE: where $W_d$ is a trainable weight matrix in the attention layer.	<a href="#">Luong 2015</a>
<b>Location-based</b>	$a_{t,i} = softmax(W_a s_t)$ *Note: This simplifies the softmax alignment to only depend on the target position.	<a href="#">Luong 2015</a>

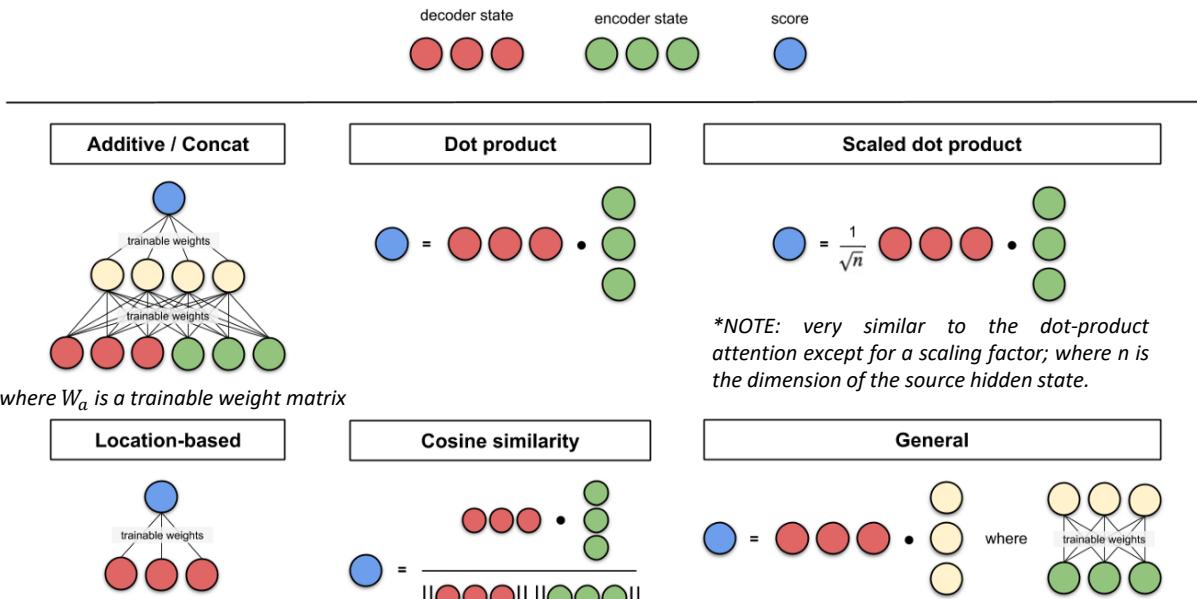
\*The papers (Luong 2015 and Vaswani 2017) can be found in the canvas content page

# 4 Attention

## Attention variants

*There are several ways to compute attention score.*

- Encoder hidden states:  $h_1, \dots, h_N \in \mathbb{R}^h$
- Decoder hidden state:  $s_t \in \mathbb{R}^h$  (on timestep  $t$ )



\*Note: This simplifies the softmax alignment to only depend on the target position.

## 4 Attention

### Categories of Attention Mechanism

*A summary of broader categories of attention mechanisms*

Name	Definition	Citation
Global or Local	<ul style="list-style-type: none"><li>• Global: Attending to the entire input state space.</li><li>• Local: Attending to the part of input state space (i.e. a patch of the input image.)</li></ul>	<b>Luong 2015</b>
Self-Attention	Relating different positions of the same input sequence. Theoretically the self-attention can adopt any attention score functions, but just replace the target sequence with the same input sequence.	<b>Cheng 2016</b>

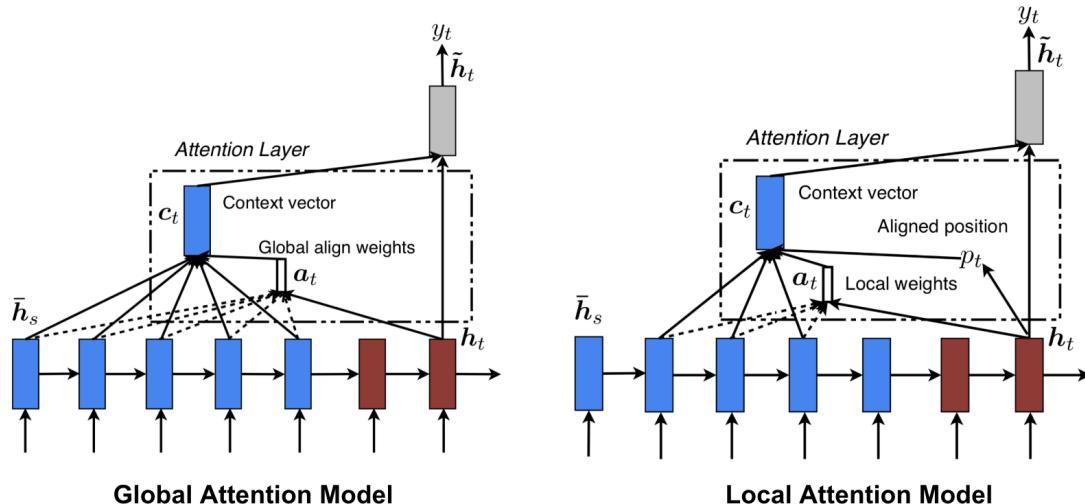
\*The papers (Luong 2015 and Cheng 2016) can be found in the canvas content page

## 4 Attention

### Categories of Attention Mechanism (1)

#### *Global/Local Attention*

- *Global: Attending to the entire input state space.*
- *Local: Attending to the part of input state space*



## 4 Attention

### Categories of Attention Mechanism (2)

#### *Self-Attention*

*The long short-term memory network (Cheng et al., 2016) paper used self-attention to do machine reading. In the example below, the self-attention mechanism enables us to learn the correlation between the current words and the previous part of the sentence.*

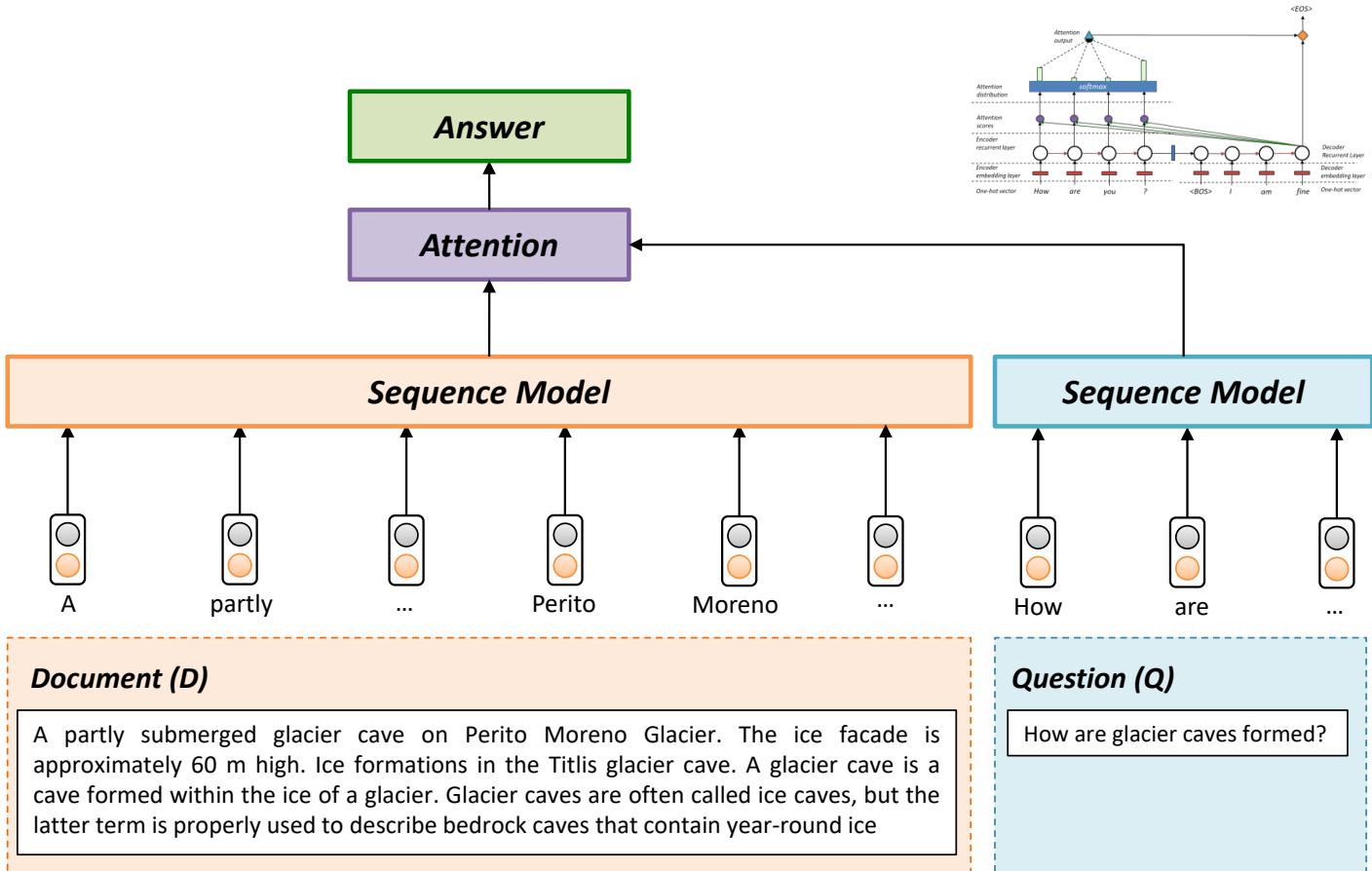
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .

*The current word is in red and the size of the blue shade indicates the activation level.*

*importance*

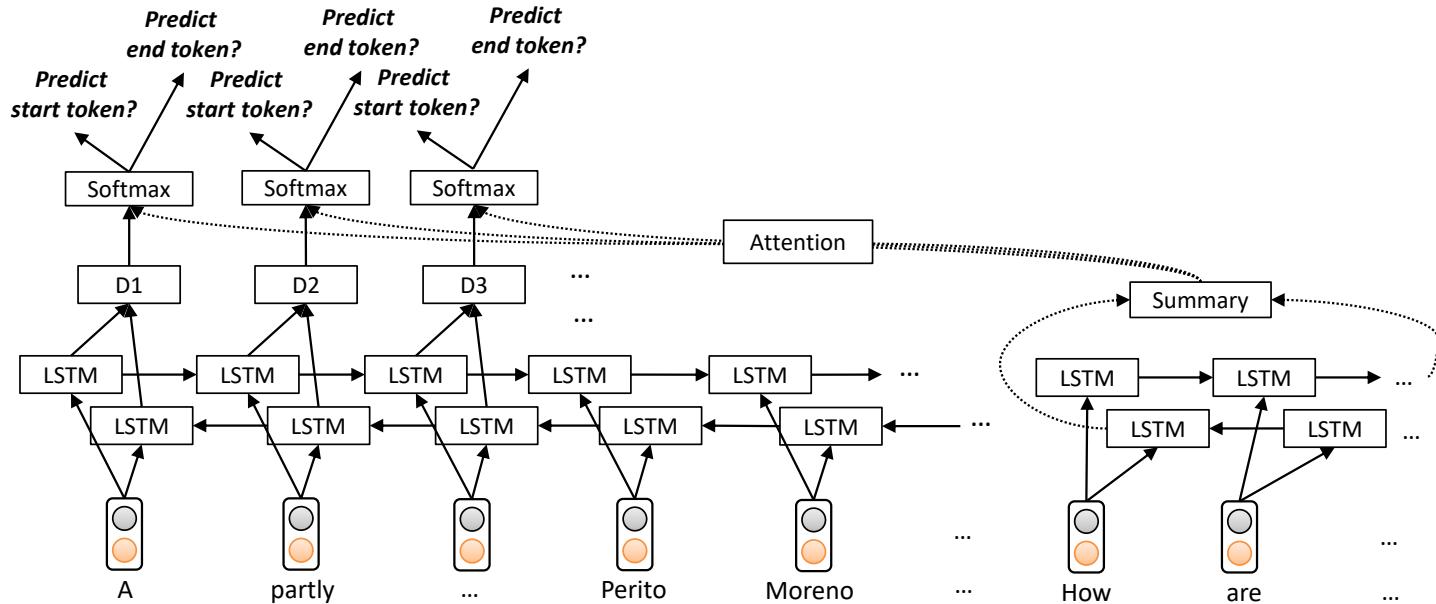
# 5 Reading Comprehension with Attention

## A Generic Neural Model for Reading Comprehension



# 5 Reading Comprehension with Attention

## Bi-LSTM for Reading Comprehension with Attention



### Document (D)

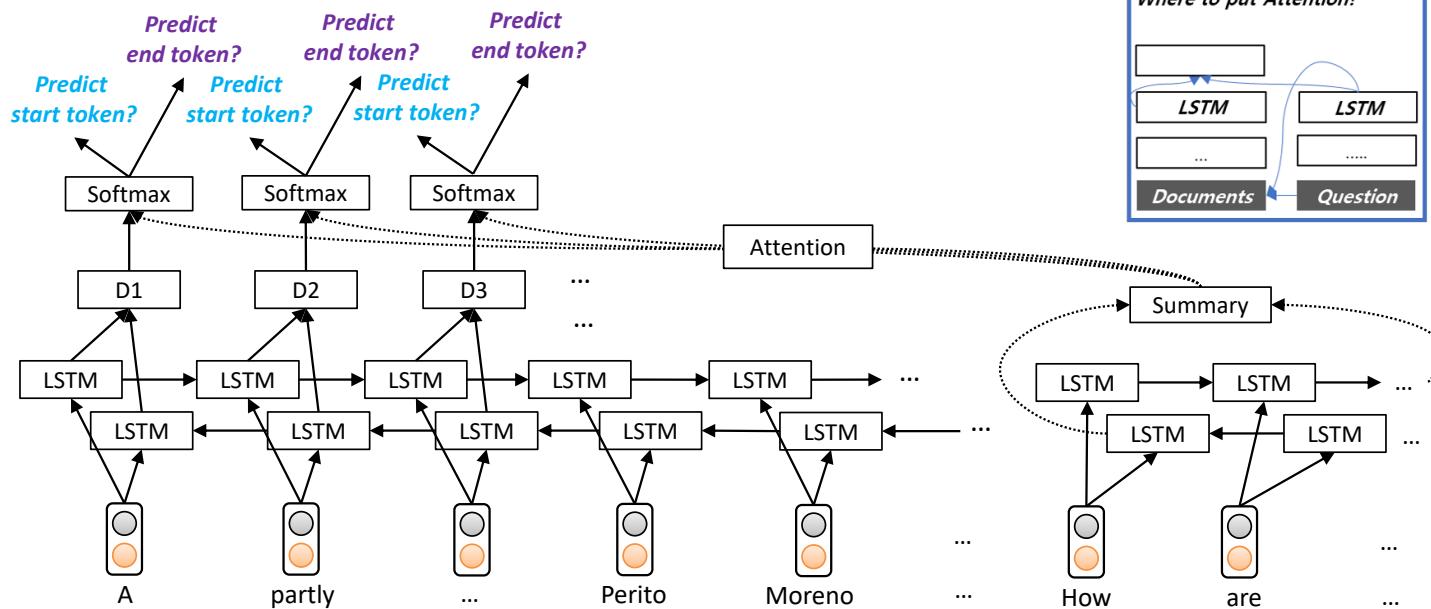
A partly submerged glacier cave on Perito Moreno Glacier. The ice facade is approximately 60 m high. Ice formations in the Titlis glacier cave. A glacier cave is a cave formed within the ice of a glacier. Glacier caves are often called ice caves, but the latter term is properly used to describe bedrock caves that contain year-round ice

### Question (Q)

How are glacier caves formed?

# 5 Reading Comprehension with Attention

## Bi-LSTM for Reading Comprehension with Attention



### Document (*D*)

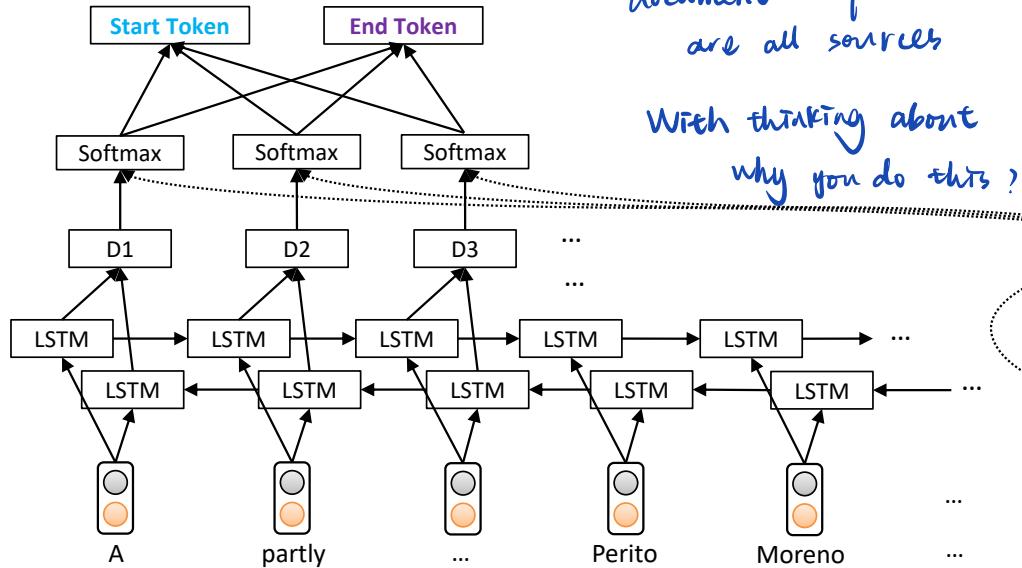
A partly submerged glacier cave on Perito Moreno Glacier. The ice facade is approximately 60 m high. Ice formations in the Titlis glacier cave. A glacier cave is a cave formed within the ice of a glacier. Glacier caves are often called ice caves, but the latter term is properly used to describe bedrock caves that contain year-round ice

### Question (*Q*)

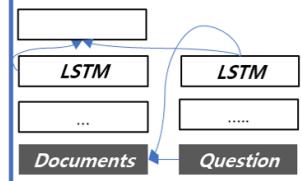
How are glacier caves formed?

# 5 Reading Comprehension with Attention

## Bi-LSTM for Reading Comprehension with Attention



Where to put Attention?



### Document (D)

A partly submerged glacier cave on Perito Moreno Glacier. The ice facade is approximately 60 m high. Ice formations in the Titlis glacier cave. A glacier cave is a cave formed within the ice of a glacier. Glacier caves are often called ice caves, but the latter term is properly used to describe bedrock caves that contain year-round ice

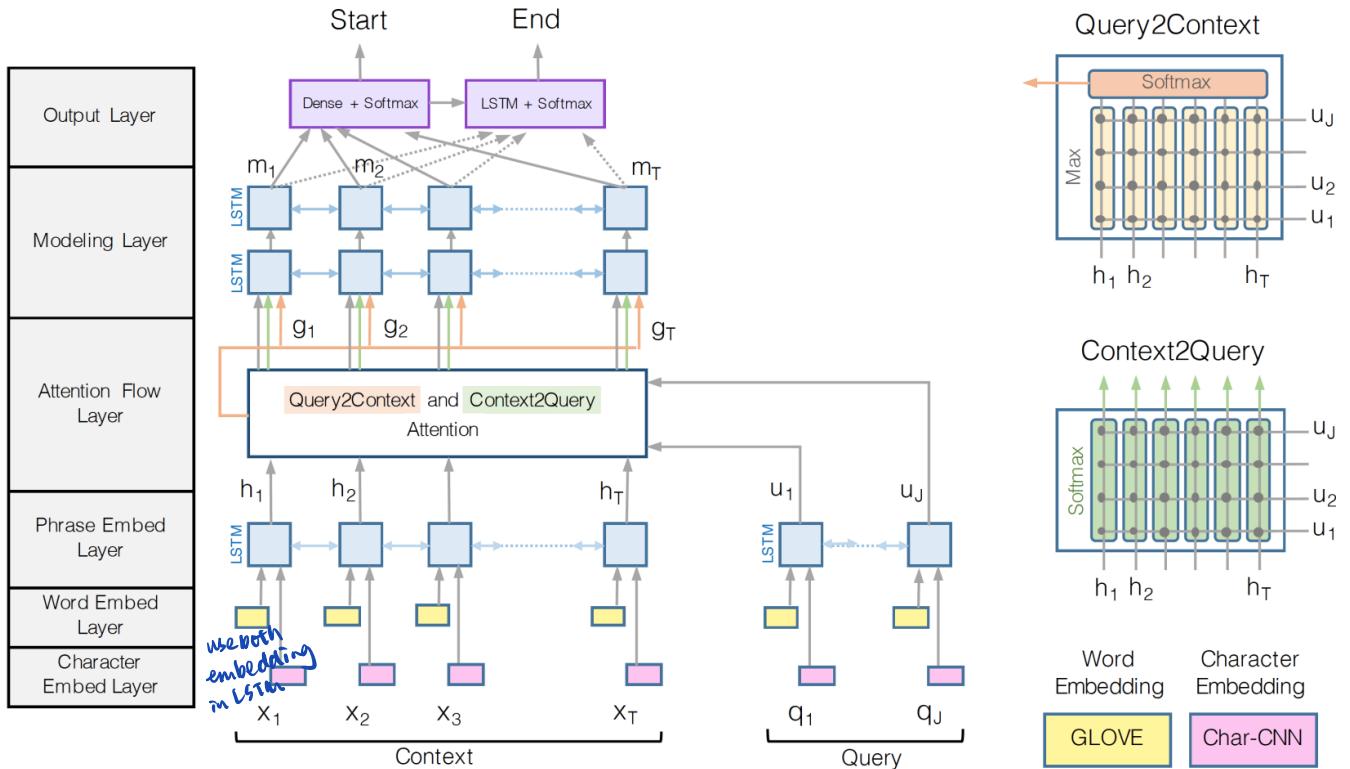
### Question (Q)

How are glacier caves formed?

# 5 Reading Comprehension with Attention

## Bi-Directional Attention Flow (Bi-DAF)

*Bi-Directional Attention Flow for Machine Comprehension (Seo et al. 2017)*



# 5 Reading Comprehension with Attention

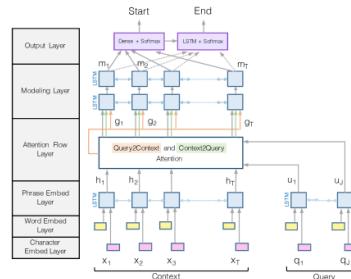
## Bi-Directional Attention Flow (Bi-DAF)

*Attention Flow layer is the core idea!*

- *Variants and improvements to the Bi-DAF architecture over the years*

Attention should flow both ways:

- 1) *the context → the question (C2Q)*
- 2) *the question → the context (Q2C)*



Both attentions are derived from a **shared similarity matrix** between the context ( $H$ ) and the query ( $U$ ), where  $S_{tj}$  indicates the similarity between t-th context word and j-th query word

$$S_{tj} = \alpha(H_{:t}, U_{:j}) \in \mathbb{R}$$

# 5 Reading Comprehension with Attention

## Bi-Directional Attention Flow (Bi-DAF)

*Attention Flow layer is the core idea!*

- *Variants and improvements to the Bi-DAF architecture over the years*

*Attention should flow both ways:*

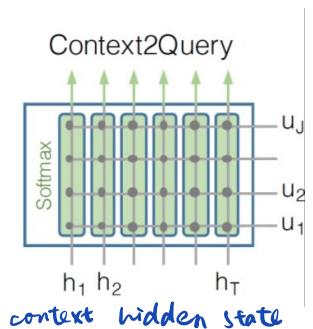
- 1) *the context → the question (C2Q)*
- 2) *the question → the context (Q2C)*

### 1. Context-to-Question (C2Q) attention:

- *which query words are most relevant to each context word*

$$\alpha^i = \text{softmax}(\mathbf{S}_{i,:}) \in \mathbb{R}^M \quad \forall i \in \{1, \dots, N\}$$

$$\mathbf{a}_i = \sum_{j=1}^M \alpha_j^i \mathbf{q}_j \in \mathbb{R}^{2h} \quad \forall i \in \{1, \dots, N\}$$



# 5 Reading Comprehension with Attention

## Bi-Directional Attention Flow (Bi-DAF)

*Attention Flow layer is the core idea!*

- *Variants and improvements to the Bi-DAF architecture over the years*

*Attention should flow both ways:*

- 1) *the context → the question (C2Q)*
- 2) *the question → the context (Q2C)*

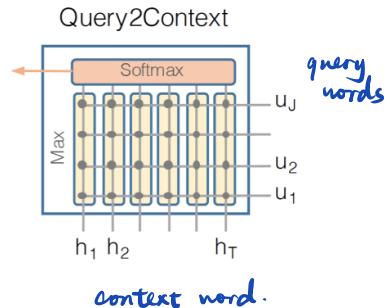
### 2. Question-to-Context (Q2C) attention:

- *the weighted sum of the most important words in the context with respect to the query – slight asymmetry through max*

$$\mathbf{m}_i = \max_j S_{ij} \in \mathbb{R} \quad \forall i \in \{1, \dots, N\}$$

$$\beta = \text{softmax}(\mathbf{m}) \in \mathbb{R}^N$$

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



# 5 Reading Comprehension with Attention

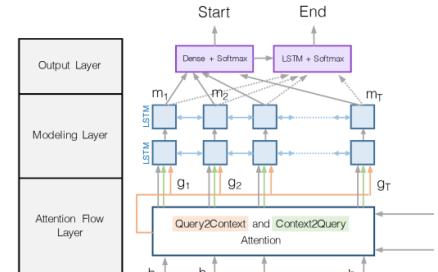
## Bi-Directional Attention Flow (Bi-DAF)

A “modelling” layer:

- Another deep (2-layer) Bi-LSTM over the passage

And answer span selection is more complex:

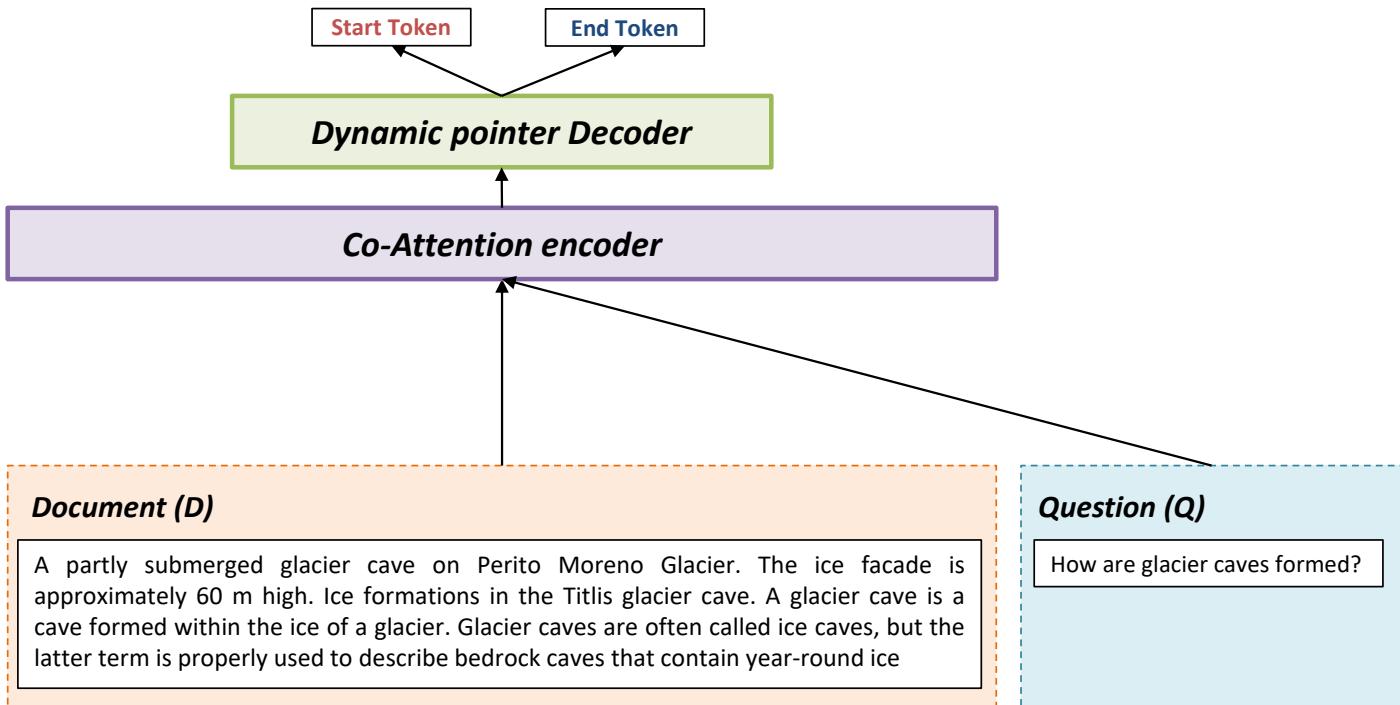
- **Start:** Pass output of BiDAF and modelling layer concatenated to a dense FF layer and then a softmax
- **End:** Put output of modelling layer  $M$  through another BiLSTM to give  $M_2$  and then concatenate with BiDAF layer and again put through dense FF layer and a softmax



## 5 Reading Comprehension with Attention

### Dynamic Coattention Networks for Question Answering (Xiong 2017)

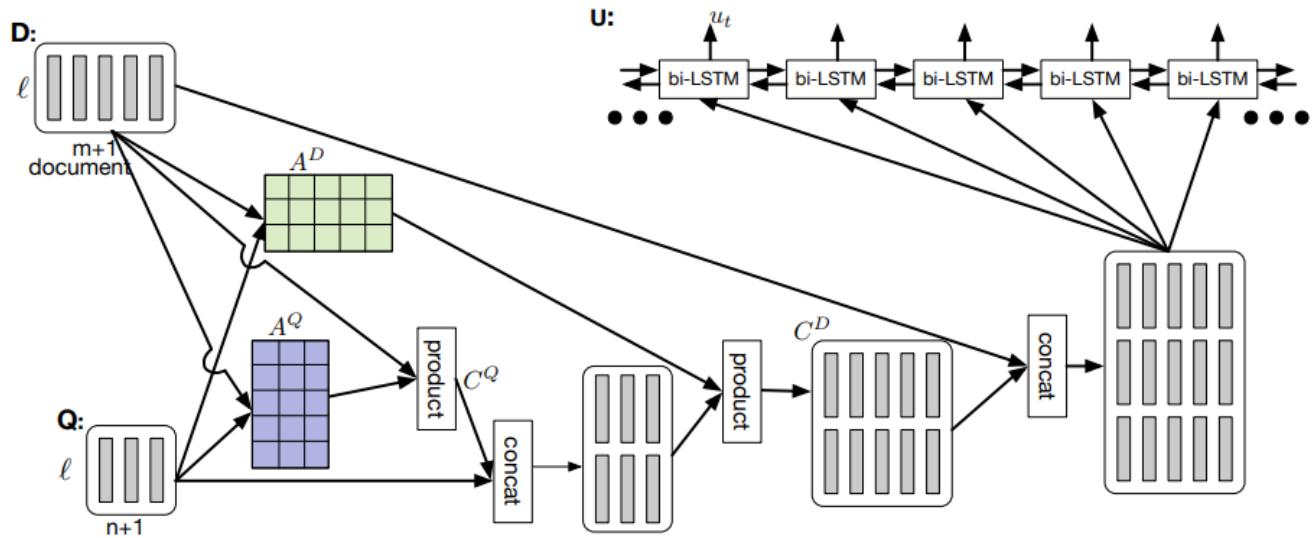
*Coattention provides a two-way attention between the context and the question.*



## 5 Reading Comprehension with Attention

### Dynamic Coattention Networks for Question Answering (Xiong 2017)

- Coattention layer again provides a two-way attention between the context and the question
- Coattention involves a second-level attention computation:
  - attending over representations that are themselves attention outputs



# / More...?

## More Advanced Architecture? Preview for following weeks

The transformer, based solely on attention mechanisms.

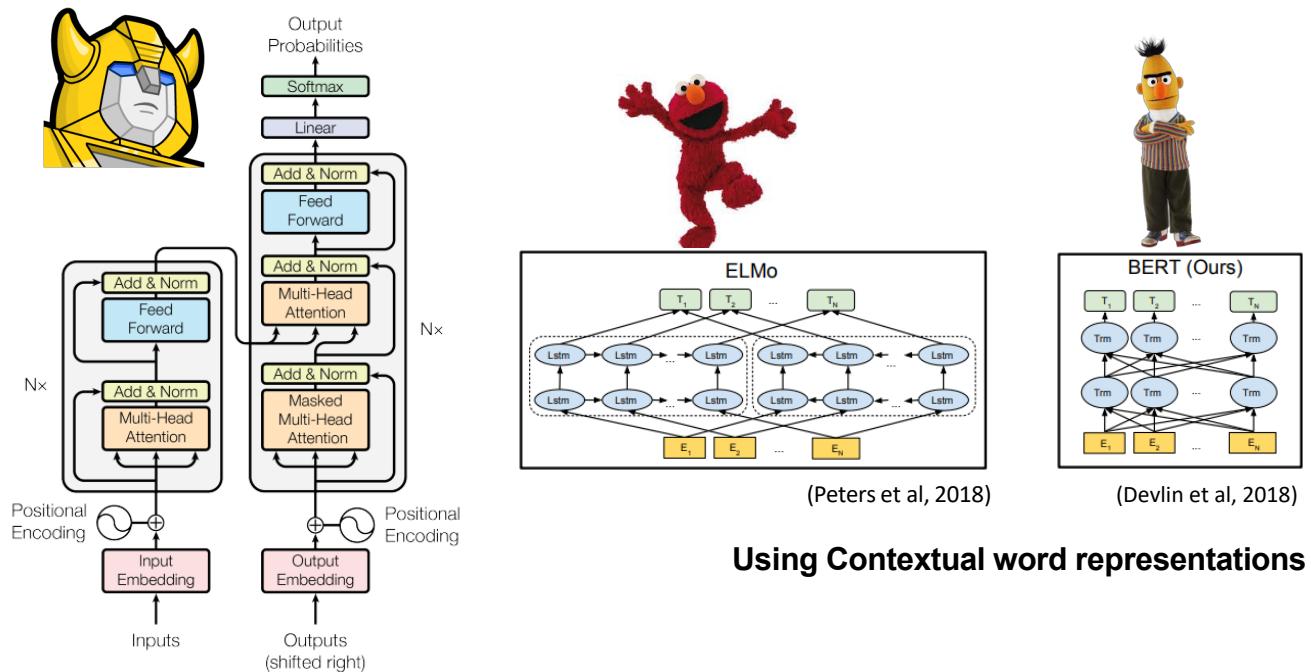


Figure 1: The Transformer - model architecture.

(Vaswani et al, 2017)

**Using Contextual word representations**

# 0 Question Answering

## Research Areas in Question Answering

Research Area	Details
Knowledge-based QA (Semantic Parsing)	<ul style="list-style-type: none"><li>• Answer is a logical form, possibly executed against a Knowledge Base</li><li>• Context is a Knowledge Base</li></ul>
Information Retrieval-based QA <ul style="list-style-type: none"><li>• Answer sentence selection</li><li>• Reading Comprehension</li></ul>	<ul style="list-style-type: none"><li>• Answer is a document, paragraph, sentence</li><li>• Context is a corpus of documents or a specific document</li></ul>
Visual QA	<ul style="list-style-type: none"><li>• <b>Answer is simple and factual</b></li><li>• <b>Context is one/multiple image(s)</b></li></ul>
Library Reference	<ul style="list-style-type: none"><li>• Answer is another question</li><li>• Context is the structured knowledge available in the library and the librarians' view of it.</li></ul>

# 6 Visual Question Answering

## Textual Question Answering: Recap

Answer questions by exploiting pure natural language.

### *Document / Passage*

*Caren watched TV last night. There was a guy playing tennis. Caren did not know who he is. He was wearing white shirts ...*

### *Question*

*What was he doing?*

### *Answer*

*Playing tennis*

# 6 Visual Question Answering

## Visual QA

Several questions require context outside of pure language.



***Question***

*What was he doing?*

***Answer***

***Playing tennis***

# 6 Visual Question Answering

## Visual QA Datasets

Recently, there are a number of visual QA datasets have sprung up. Some of the more popular ones include:

### Type #1: Real images

VQA v2.0 (Goyal et al. 2017)

Where is the child sitting?  
fridge arms



How many children are in the bed?

2



1



### Type #2: Semantic reasoning

CLEVER (Johnson et al. 2016)



Q: Are there an **equal number** of **large things** and **metal spheres**?

Q: **What size** is the **cylinder** that is **left of** the **brown metal thing** that is **left of** the **big sphere**?

Q: There is a **sphere** with the **same size** as the **metal cube**; is it **made of the same material** as the **small red sphere**?

### Type #3: Combined

GQA (Hudson and Manning, 2019)



Is the **bowl** to the right of the **green apple**?

What type of **fruit** in the image is **round**?

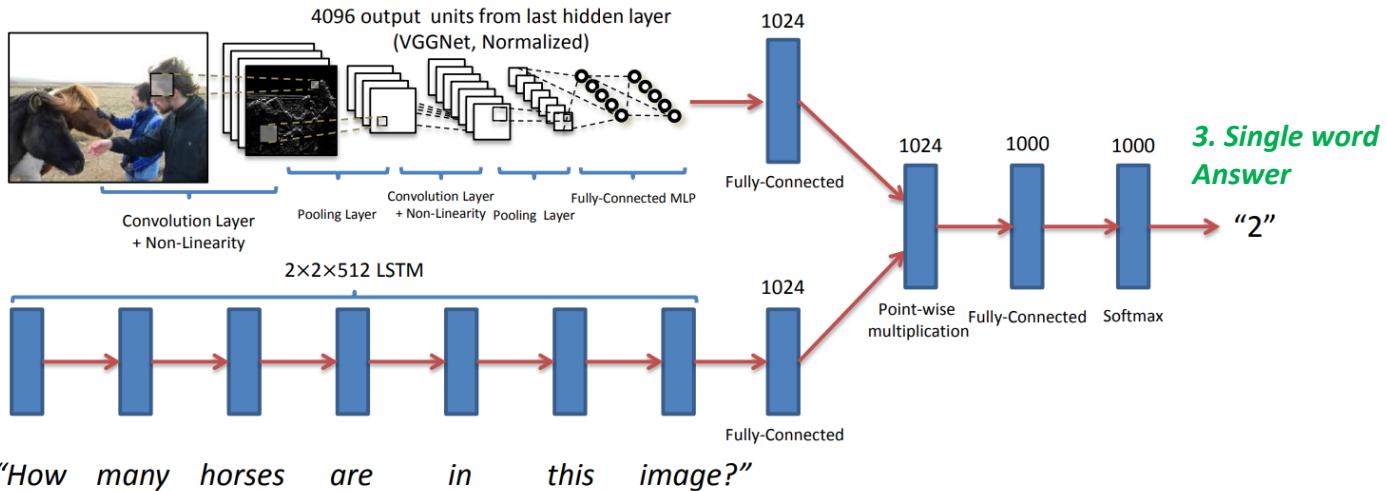
What color is the **fruit** on the right side, **red** or **green**?

Is there any **milk** in the **bowl** to the left of the **apple**?

# 6 Visual Question Answering

How does it work?

**2. Context is a single picture using convolutional neural network (CNN)**

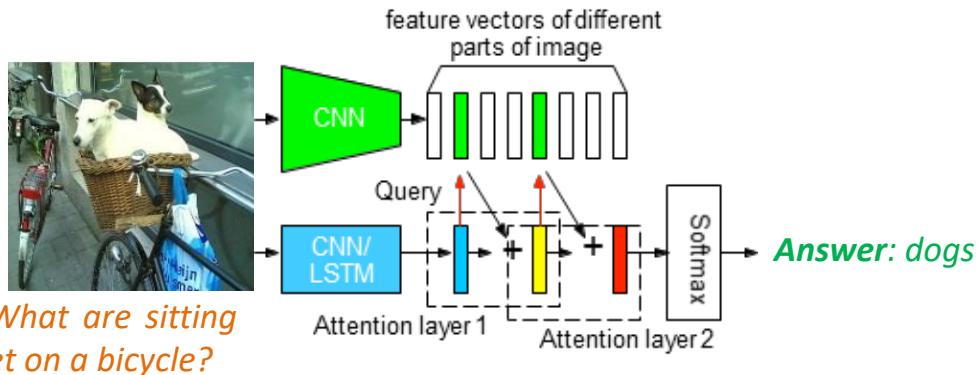


**1. Encode sentence with sequence models**

# 6 Visual Question Answering

## How does it work?

*The idea of Visual QA is exactly same as reading comprehension-oriented.  
Why can't we use **Attention** then? (Yang et al. 2015)*



# 6 Visual Question Answering

## Visual QA with Attention

Let's try some example. VisualQA (<http://vqa.cloudcv.org/>)

- ( a ) What are pulling a man on a wagon down on dirt road?  
 Answer: horses Prediction: horses



- ( b ) What is the color of the box ?  
 Answer: red Prediction: red



- ( c ) What next to the large umbrella attached to a table?  
 Answer: trees Prediction: tree



- ( d ) How many people are going up the mountain with walking sticks?  
 Answer: four Prediction: four



- ( e ) What is sitting on the handle bar of a bicycle?  
 Answer: bird Prediction: bird



- ( f ) What is the color of the horns?  
 Answer: red Prediction: red



Original Image

First Attention Layer

Second Attention Layer

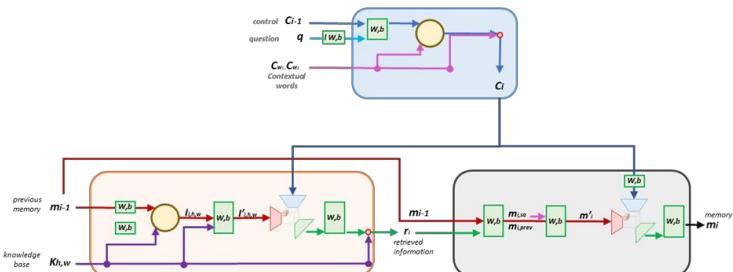
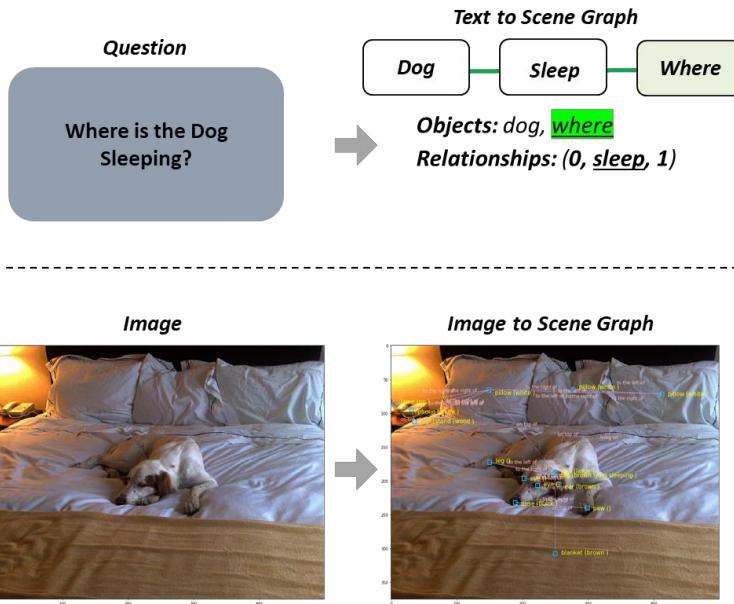
Original Image

First Attention Layer

Second Attention Layer

# 6 Visual Question Answering

Visual QA with Attention (Usyd NLP Group 2020)



# 6 Visual Question Answering

## Visual QA with Attention (Usyd NLP Group 2020)

Predictions #: 44  
 Image ID: 2356967  
 Object list index #: 62642  
 Question: What color is that sky?  
 Prediction: gray  
 Answer: gray

**Step 1**



**Step 2**



**Step 3**



**Step 4**



What	color	is	that	sky
------	-------	----	------	-----

What	color	is	that	sky
------	-------	----	------	-----

What	color	is	that	sky
------	-------	----	------	-----

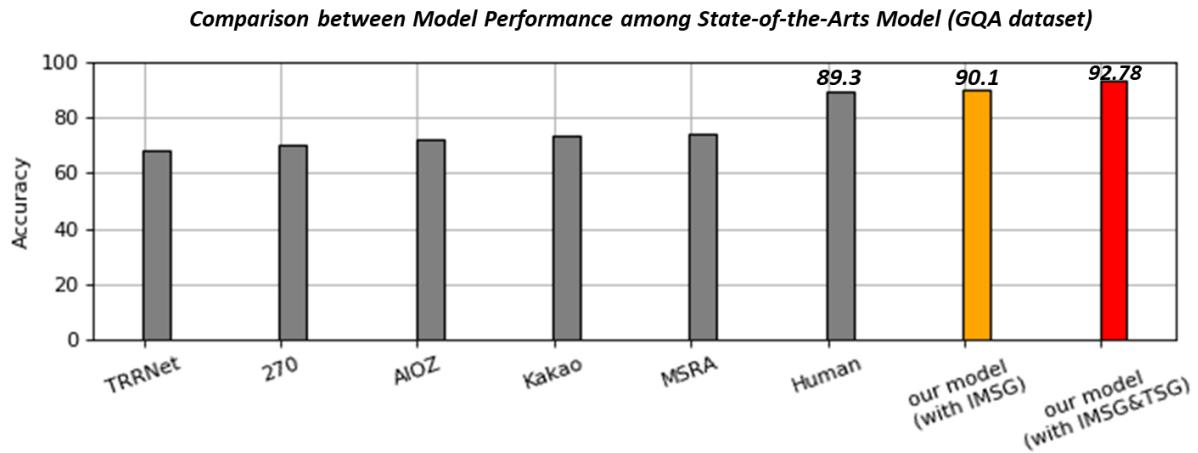
What	color	is	that	sky
------	-------	----	------	-----

Let's have a look at the demo!!

# 6 Visual Question Answering

## Visual QA with Attention (Usyd NLP Group 2020)

### *Testing Result*



# / Additional QA

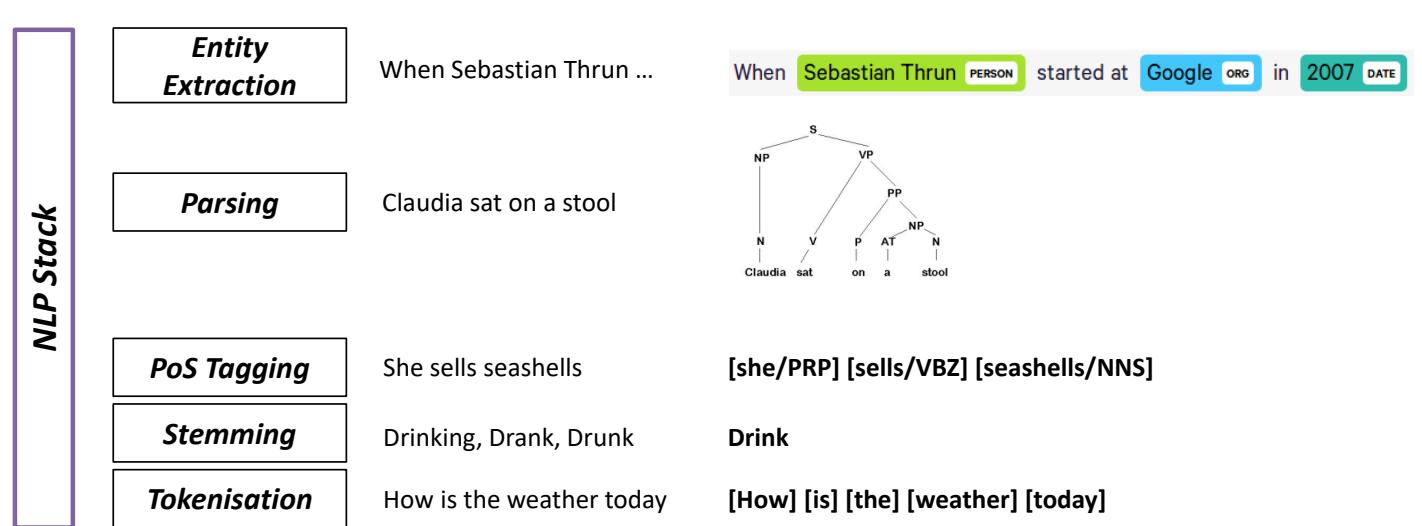
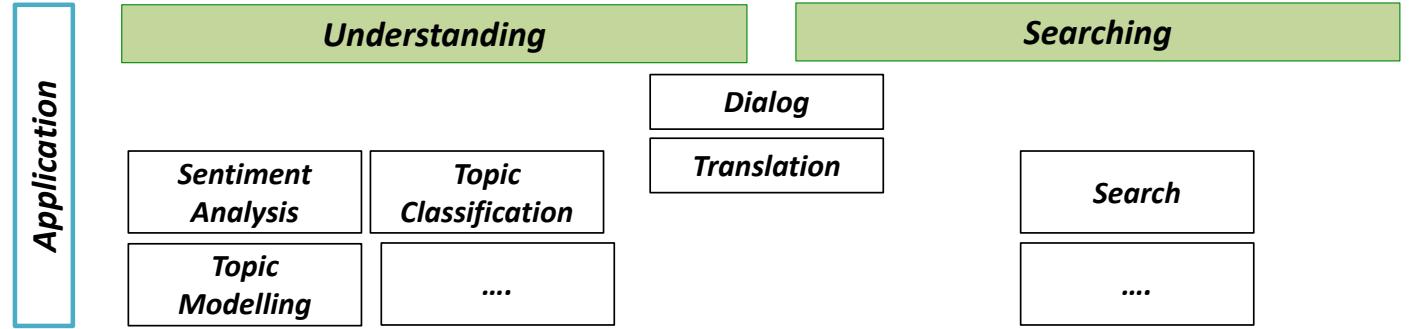
**There is no reason to limit to just IR-based or Knowledge-based**

*Using multiple information sources? IBM Watson!*



# / The big picture of NLP

## The purpose of Natural Language Processing: Overview



# / Reference

## Reference for this lecture

- Deng, L., & Liu, Y. (Eds.). (2018). Deep Learning in Natural Language Processing. Springer.
- Rao, D., & McMahan, B. (2019). Natural Language Processing with PyTorch: Build Intelligent Language Applications Using Deep Learning. " O'Reilly Media, Inc.".
- Manning, C. D., Manning, C. D., & Schütze, H. (1999). Foundations of statistical natural language processing. MIT press.
- Manning, C 2018, Natural Language Processing with Deep Learning, lecture notes, Stanford University
  
- Berant, J., Chou, A., Frostig, R., & Liang, P. (2013). Semantic parsing on freebase from question-answer pairs. In Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing (pp. 1533-1544).
- Chen, D., Bolton, J., & Manning, C. D. (2016). A thorough examination of the cnn/daily mail reading comprehension task. arXiv preprint arXiv:1606.02858.
- Chen, D., Fisch, A., Weston, J., & Bordes, A. (2017). Reading wikipedia to answer open-domain questions. arXiv preprint arXiv:1704.00051.
- Seo, M., Kembhavi, A., Farhadi, A., & Hajishirzi, H. (2016). Bidirectional attention flow for machine comprehension. arXiv preprint arXiv:1611.01603.
- Yang, Z., He, X., Gao, J., Deng, L., & Smola, A. (2016). Stacked attention networks for image question answering. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 21-29).
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In Advances in neural information processing systems (pp. 5998-6008).
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.