# 1-of-N Encoding

apple = [ 1  0  0  0  0]
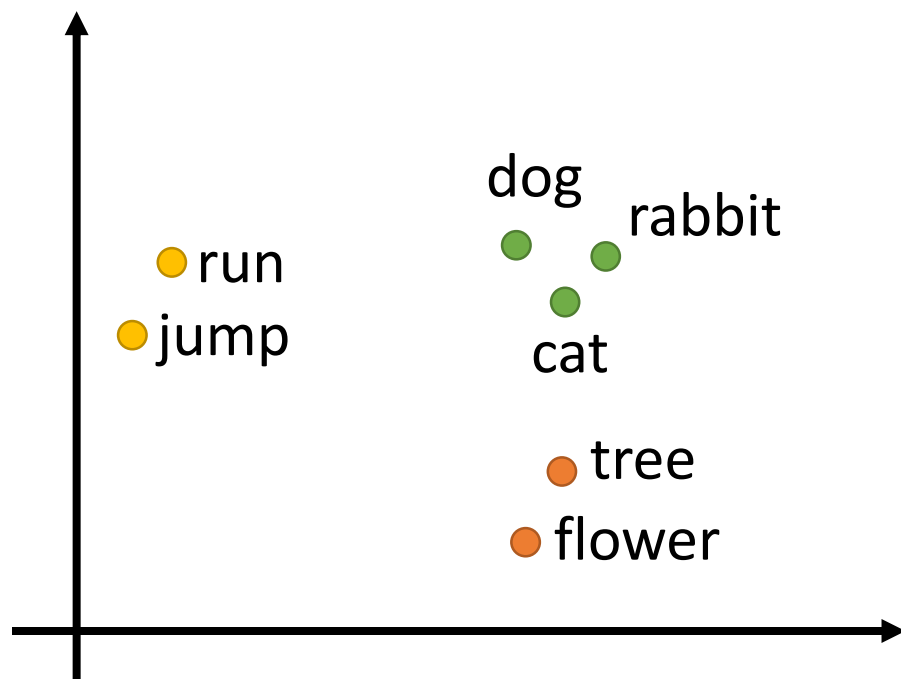
bag = [ 0  1  0  0  0]

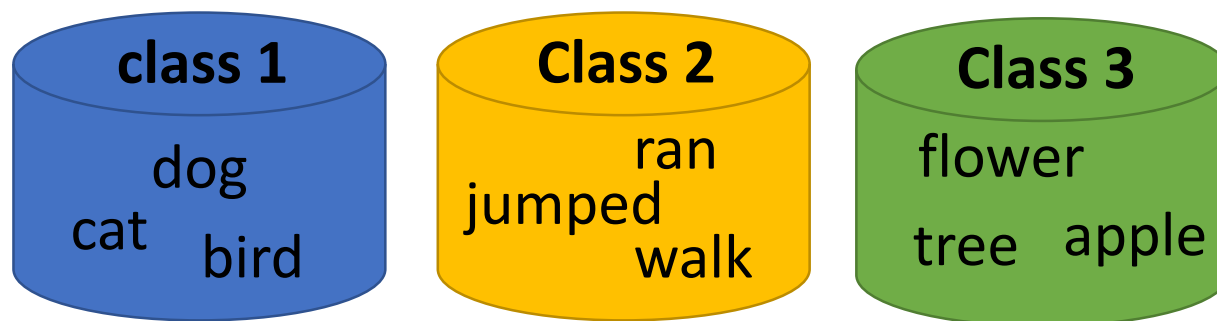cat = [ 0  0  1  0  0]

dog = [ 0  0  0  1  0]

elephant = [ 0  0  0  0  1]

# Word Embedding

run

jump

dog

rabbit

cat

tree

flower

# Word Class

**class 1**

dog
cat
bird

**Class 2**

ran
jumped
walk

**Class 3**

flower
tree  apple

# A word can have multiple senses.

Have you paid that money to the bank yet ?
It is safest to deposit your money in the bank .

The victim was found lying dead on the river bank .
They stood on the river bank to fish.

The hospital has its own blood bank.

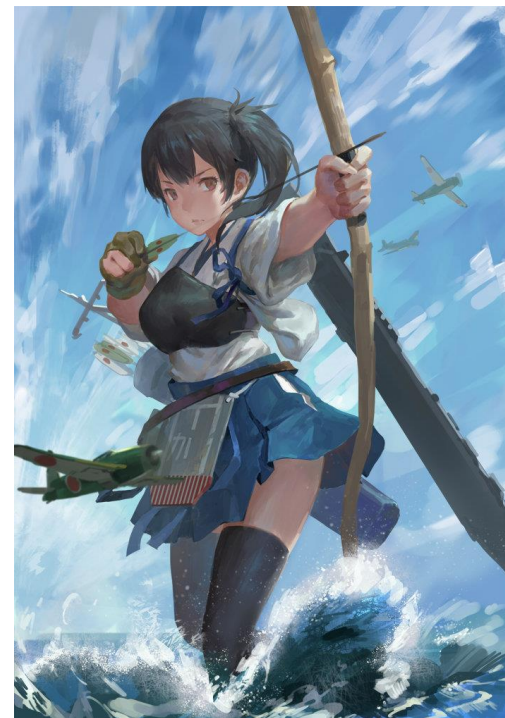The third sense or not?

# More Examples





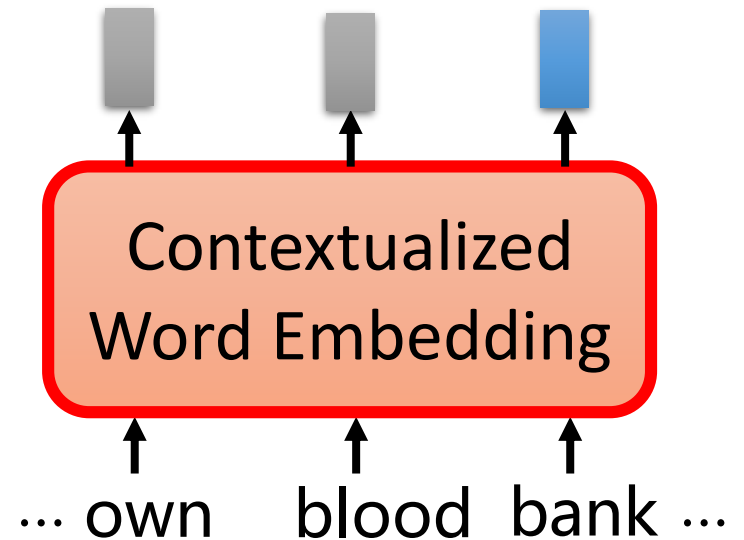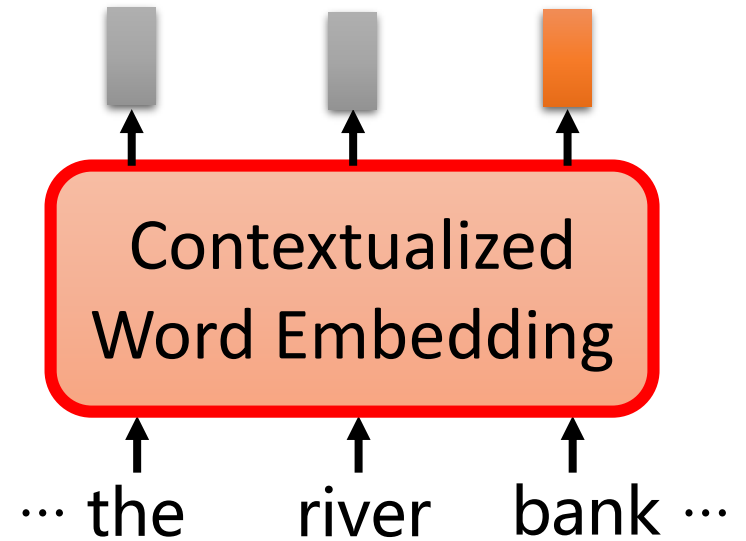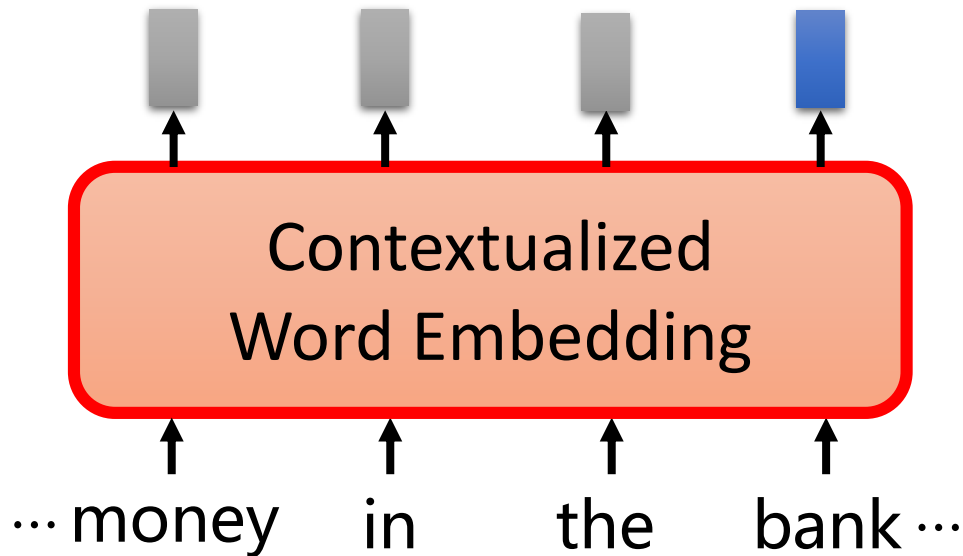這是
加賀號護衛艦

他是尼祿



她也是尼祿

這也是加賀
號護衛艦

# Contextualized Word Embedding

# Embeddings from Language Model (ELMO)

https://arxiv.org/abs/1802.05365

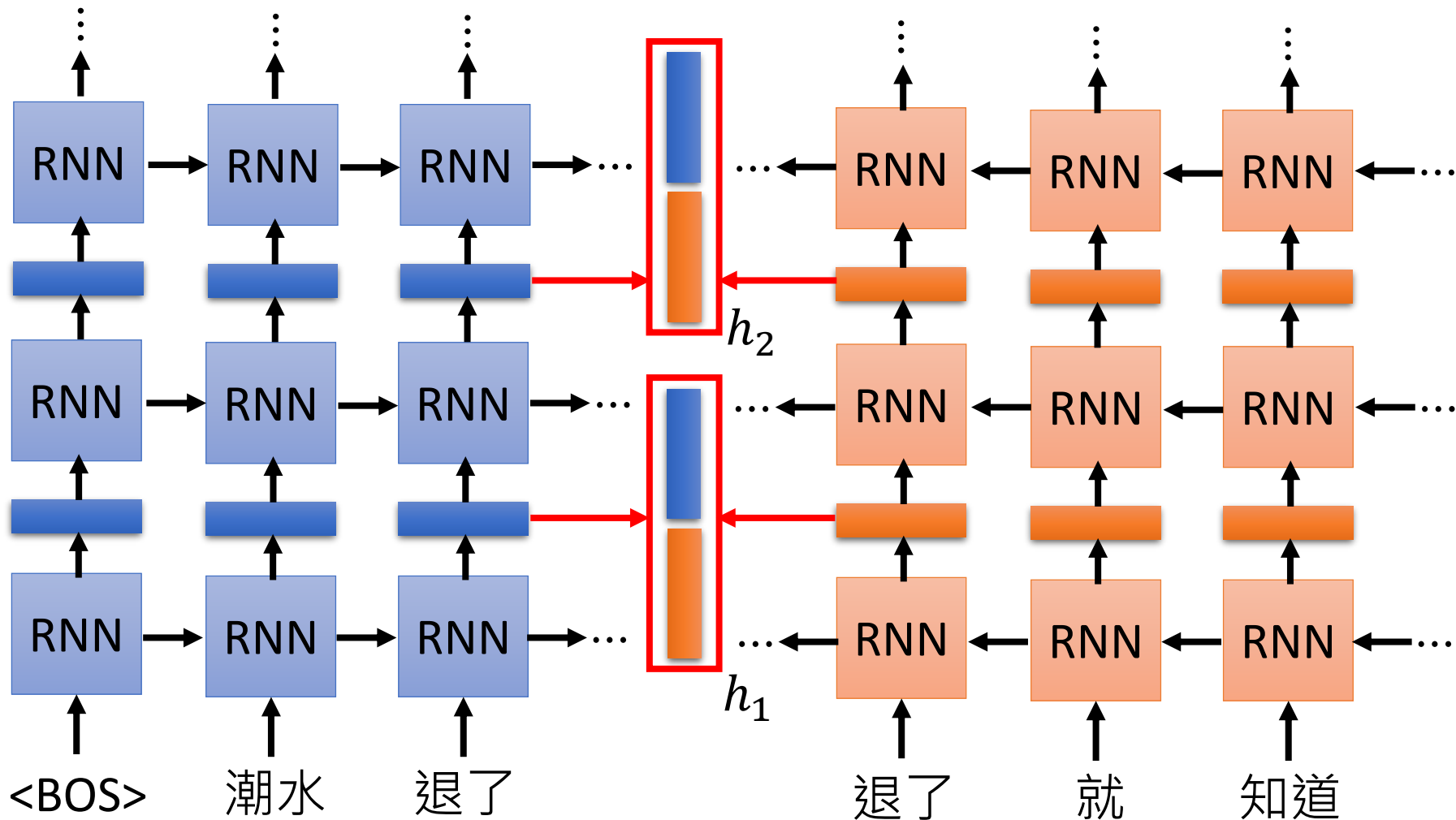- RNN-based language models (trained from lots of sentences)

    e.g. given "潮水 退了 就 知道 誰 沒穿 褲子"

# ELMO

Each layer in deep LSTM can generate a latent representation.

Which one should we use???

我全都要

# ELMO

$$\blacksquare = \alpha_1 \blacksquare + \alpha_2 \blacksquare$$

Learned with the down stream tasks



ELMO

潮水　退了　就　知道 ……



small                                    large

# Bidirectional Encoder Representations from Transformers (BERT)

- BERT = Encoder of Transformer

Learned from a large amount of text without annotation

# Training of BERT

- Approach 1:
  Masked LM

# _Training of BERT_

## Approach 2: Next Sentence Prediction

yes

Linear Binary Classifier

[CLS]: the position that outputs classification results

[SEP]: the boundary of two sentences

Approaches 1 and 2 are used at the same time.

BERT

[CLS]　醒醒　吧　[SEP]　你　沒有　妹妹

# *Training of BERT*

## Approach 2: Next Sentence Prediction

No

[CLS]: the position that outputs classification results

[SEP]: the boundary of two sentences

Linear Binary Classifier

Approaches 1 and 2 are used at the same time.

BERT

[CLS]　醒醒　吧　[SEP]　眼睛　業障　重

# How to use BERT – Case 1



class

Linear Classifier

Trained from Scratch

BERT ➤ Fine-tune

[CLS]  $W_1$  $W_2$  $W_3$

sentence

Input: single sentence, output: class

Example:
Sentiment analysis (our HW),
Document Classification

# How to use BERT – Case 2



class        class        class

Linear Cls   Linear Cls   Linear Cls

BERT

[CLS]   $w_1$   $w_2$   $w_3$

sentence

Input: single sentence, output: class of each word

Example: Slot filling

arrive   Taipei   on   November   2nd

other   dest   other   time   time

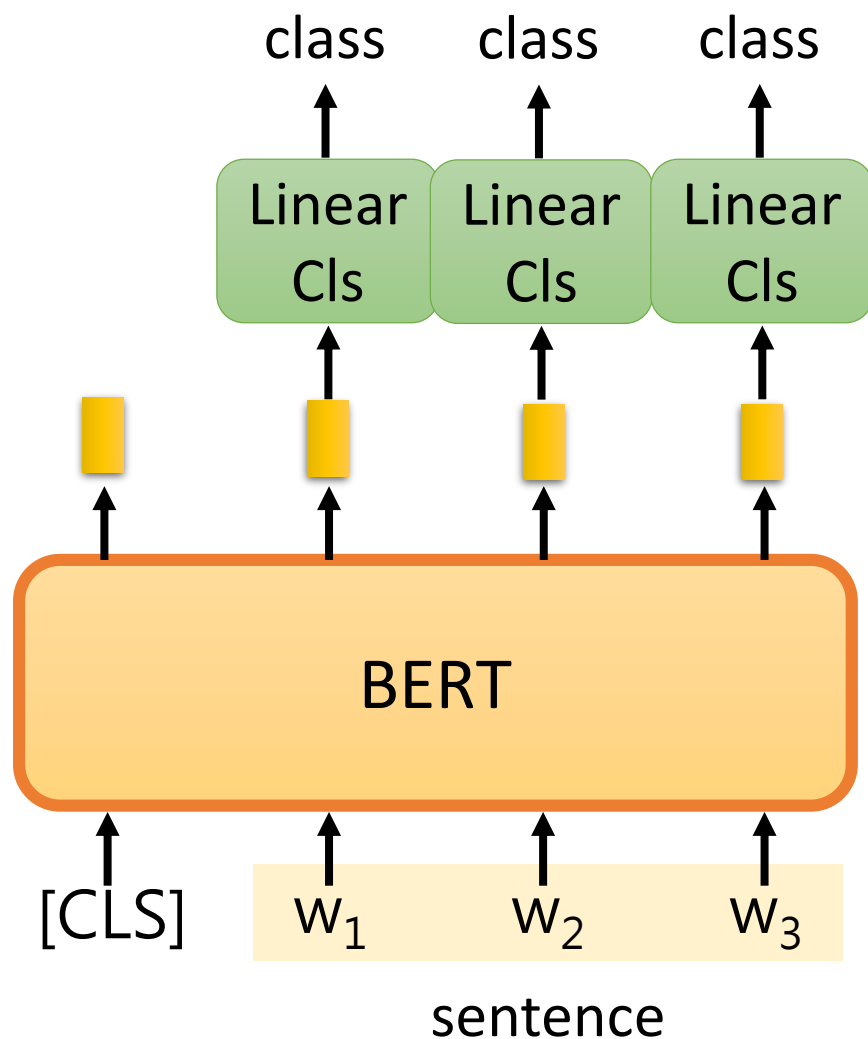# How to use BERT – Case 3



Class

Linear Classifier

Input: two sentences, output: class

Example: Natural Language Inference

Given a "premise", determining whether a "hypothesis" is T/F/ unknown.

BERT

[CLS]    $W_1$    $W_2$    [SEP]    $W_3$    $W_4$    $W_5$

Sentence 1          Sentence 2

# How to use BERT – Case 4

- Extraction-based Question Answering (QA) (E.g. SQuAD)

**Document**: $D = \{d_1, d_2, \cdots, d_N\}$

**Query**: $Q = \{q_1, q_2, \cdots, q_N\}$
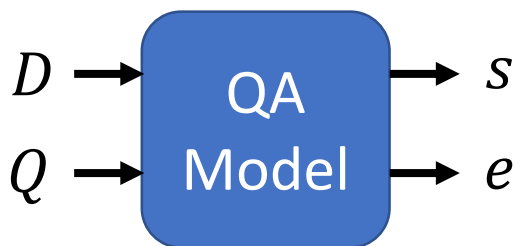
$D \rightarrow$ [QA Model] $\rightarrow s$

$Q \rightarrow$ [QA Model] $\rightarrow e$

output: two integers $(s, e)$

**Answer**: $A = \{q_s, \cdots, q_e\}$

In meteorology, precipitation is any product of the condensation of [17] spheric 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 [77] atte [79] cations are called "showers".

What causes precipitation to fall?
**gravity**     $s = 17, e = 17$

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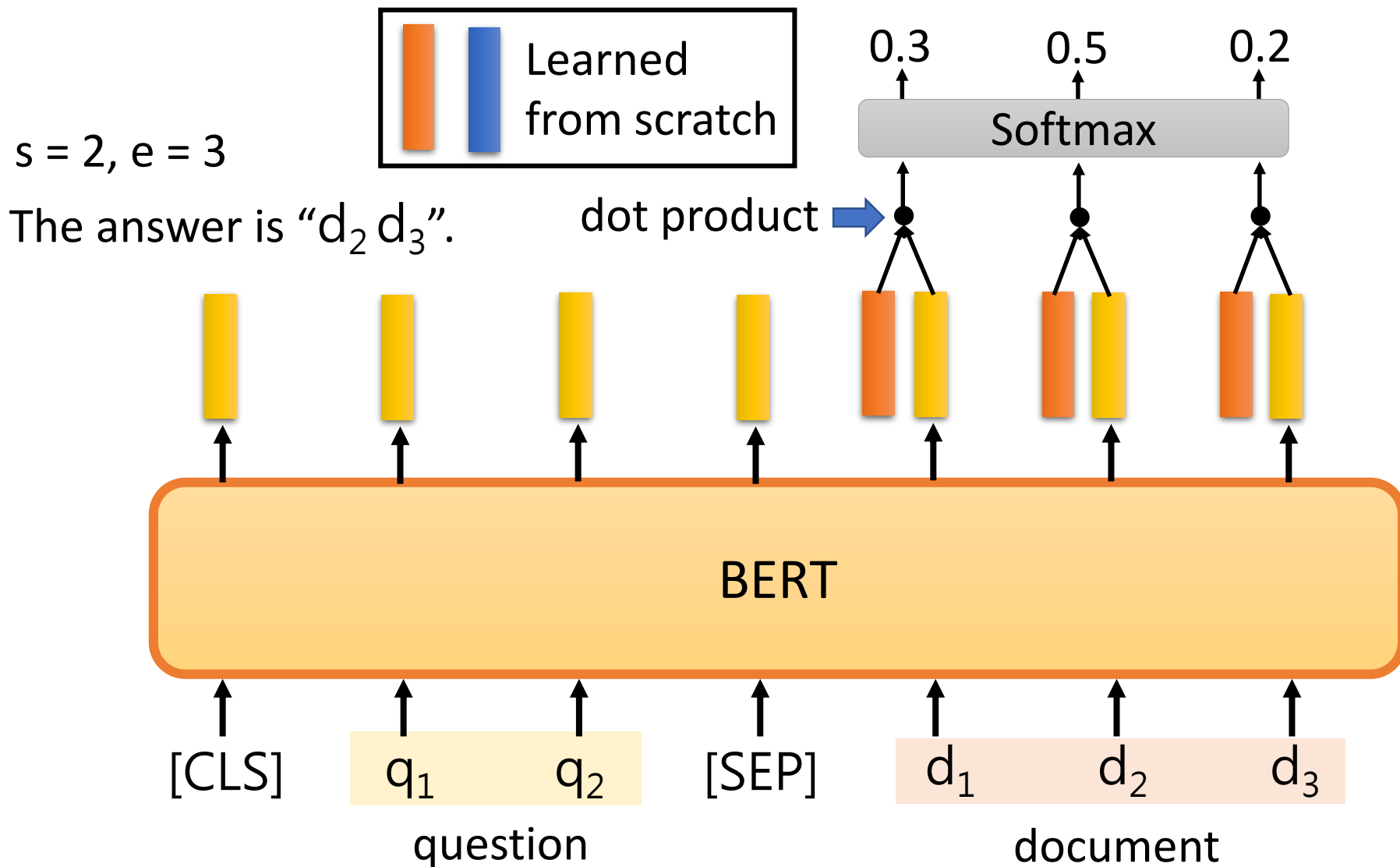**graupel**

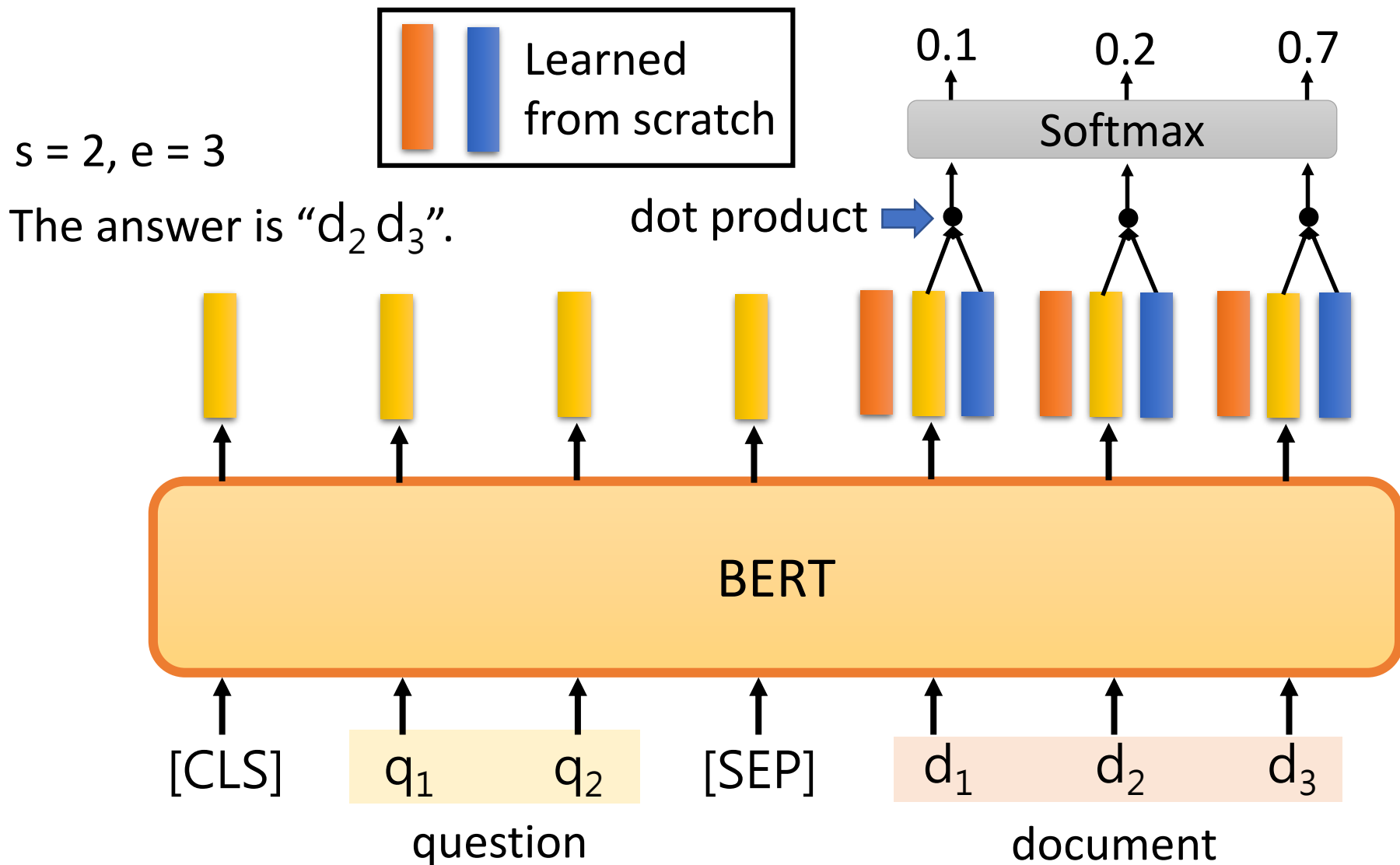Where do water droplets collide with ice crystals to form precipitation?
**within a cloud**     $s = 77, e = 79$
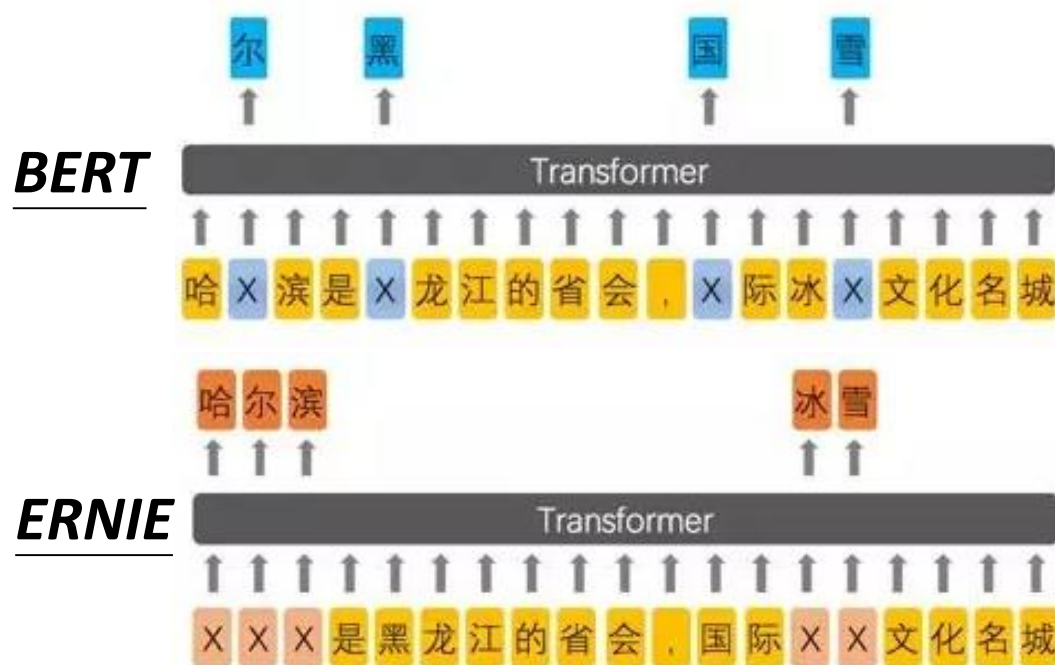
# How to use BERT – Case 4

s = 2, e = 3

The answer is "$d_2$ $d_3$".

# How to use BERT – Case 4

# BERT 屠榜 ……

| Rank | Model | EM | F1 |
|---|---|---|---|
| | Human Performance<br>*Stanford University*<br>(Rajpurkar & Jia et al. '18) | 86.831 | 89.452 |
| 1<br>Mar 20, 2019 | BERT + DAE + AoA (ensemble)<br>*Joint Laboratory of HIT and iFLYTEK Research* | **87.147** | **89.474** |
| 2<br>Mar 15, 2019 | BERT + ConvLSTM + MTL + Verifier (ensemble)<br>*Layer 6 AI* | 86.730 | 89.286 |
| 3<br>Mar 05, 2019 | BERT + N-Gram Masking + Synthetic Self-<br>Training (ensemble)<br>*Google AI Language*<br>https://github.com/google-research/bert | 86.673 | 89.147 |
| 4<br>May 21, 2019 | XLNet (single model)<br>*XLNet Team* | 86.346 | 89.133 |
| 5<br>Apr 13, 2019 | SemBERT(ensemble)<br>*Shanghai Jiao Tong University* | 86.166 | 88.886 |

SQuAD 2.0

# Enhanced Representation through Knowledge Integration (ERNIE)
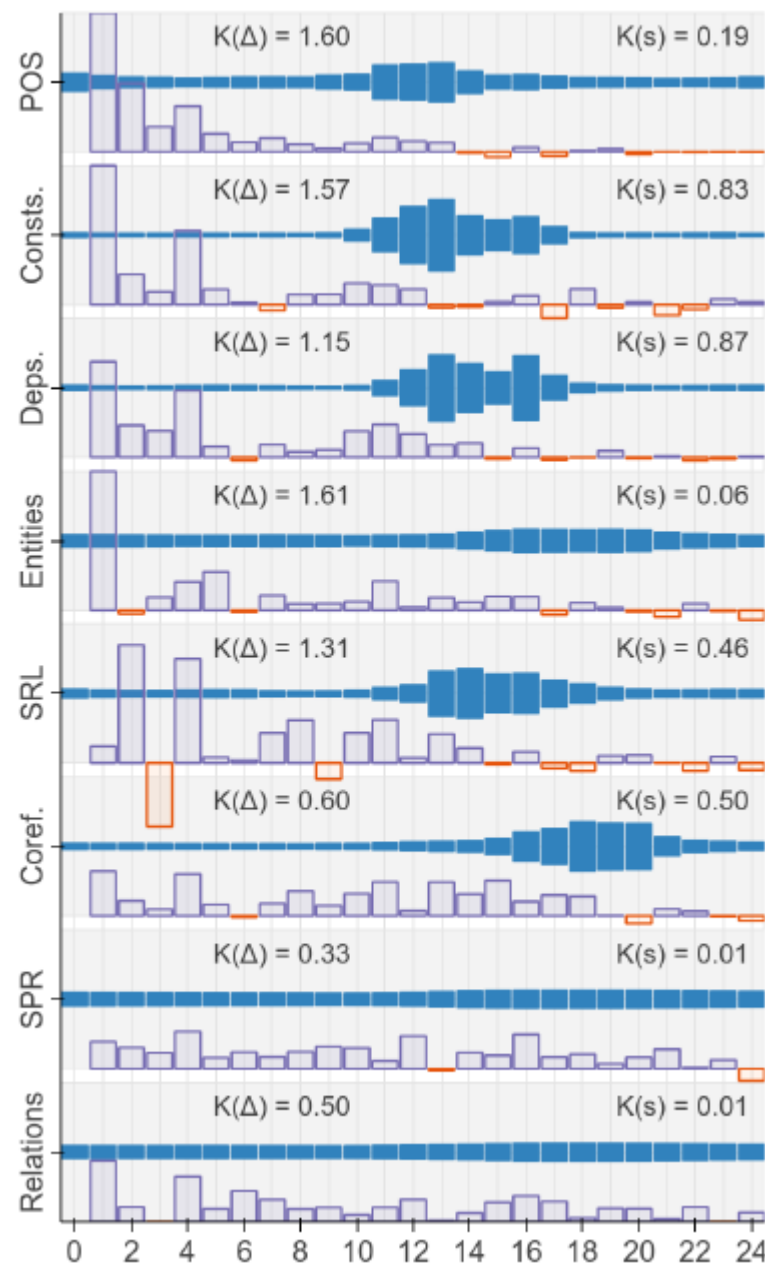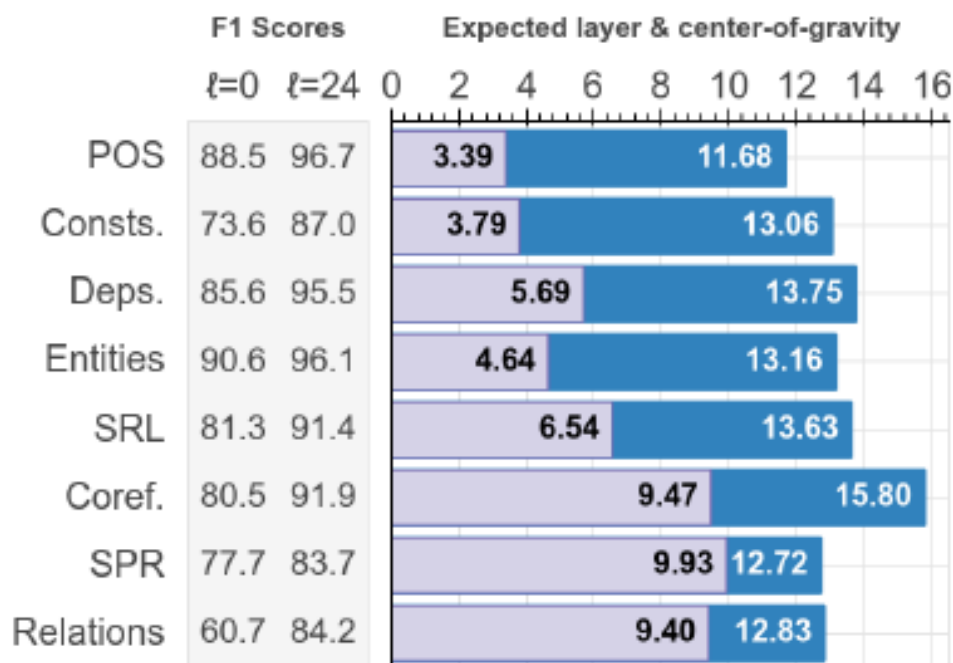
- Designed for Chinese



**BERT**

**ERNIE**

Source of image:
https://zhuanlan.zhihu.com/p/59436589

https://arxiv.org/abs/1904.09223

# What does BERT learn?

https://arxiv.org/abs/1905.05950
https://openreview.net/pdf?id=SJzSgnRcKX

# Multilingual BERT

Trained on 104 languages



Task specific training data for English

En Class 1

En Class 2

En Class 3

Task specific testing data for Chinese

Zh ?

Zh ?

Zh ?

Zh ?

Zh ?

# Generative Pre-Training (GPT)



Transformer Decoder

BERT
(340M)

ELMO
(94M)

GPT-2
(1542M)

Source of image: https://huaban.com/pins/1714071707/

# Generative Pre-Training (GPT)

退了

Many Layers ...

# Generative Pre-Training (GPT)

就

Many Layers ...

$b^3$

$\hat{\alpha}_{3,1}$  $\hat{\alpha}_{3,2}$  $\hat{\alpha}_{3,3}$

$\times$  $\times$  $\times$

$q^1$ $k^1$ $v^1$    $q^2$ $k^2$ $v^2$    $q^3$ $k^3$ $v^3$    $q^4$ $k^4$ $v^4$
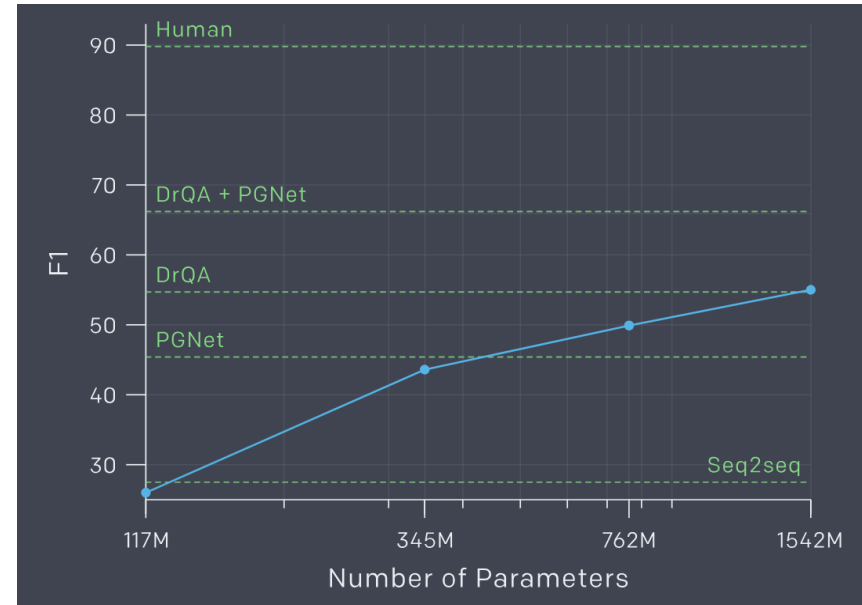
$a^1$    $a^2$    $a^3$    $a^4$

\<BOS>    潮水    退了    就

# Zero-shot Learning?

- ***Reading Comprehension***

$d_1, d_2, \cdots, d_N,$
"Q:", $q_1, q_2, \cdots, q_N,$
"A:"



CoQA

- ***Summarization***   $d_1, d_2, \cdots, d_N,$"TL;DR:"

- ***Translation***

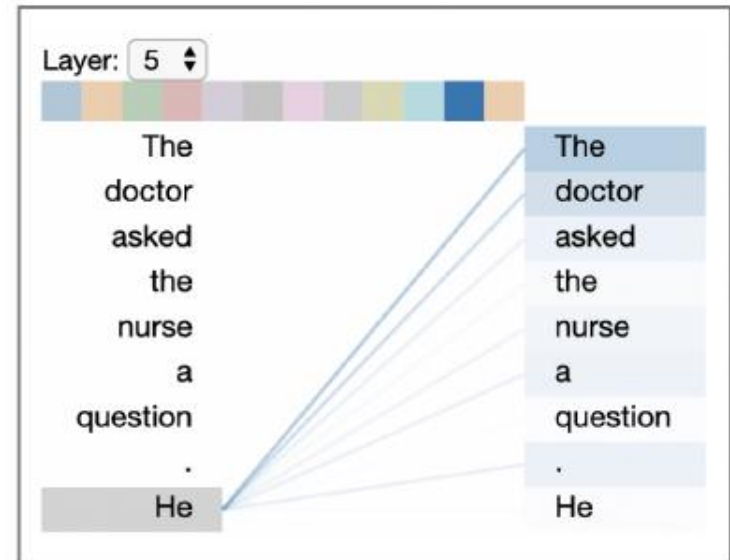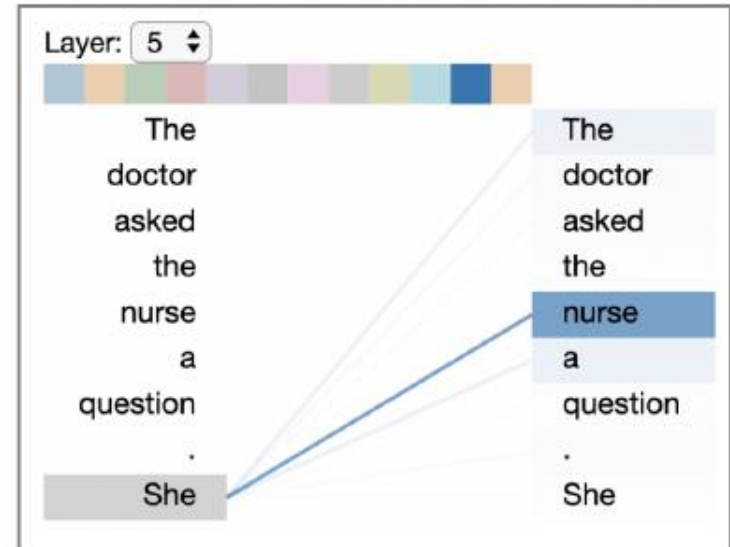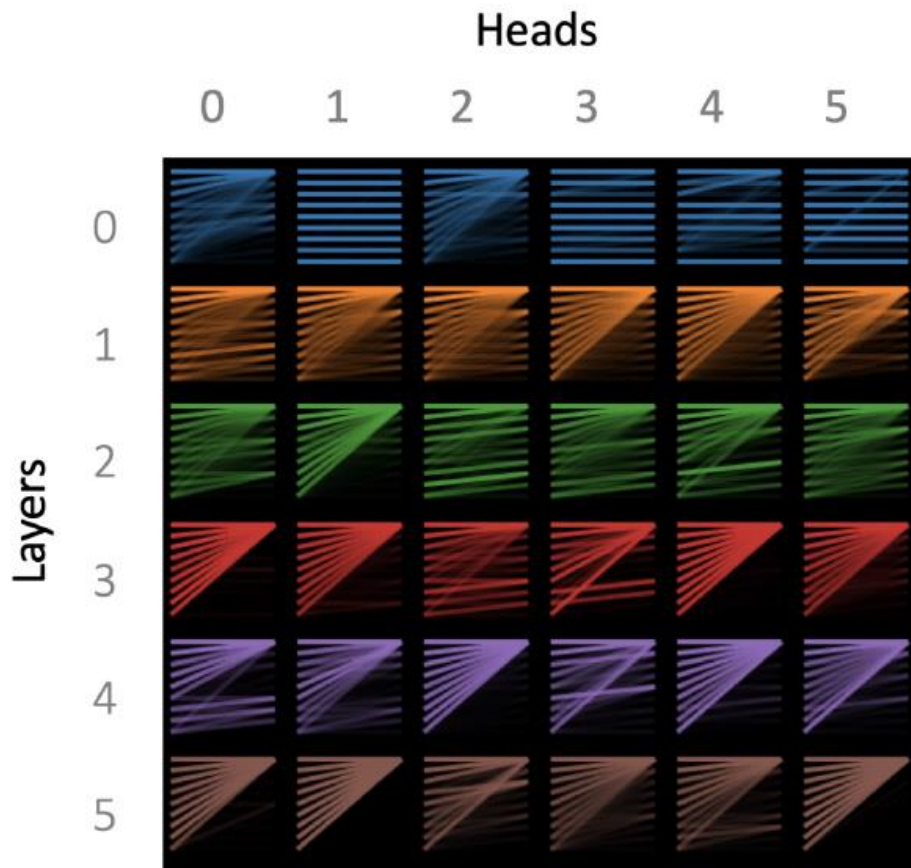| English sentence 1 | = | French sentence 1 |
| English sentence 2 | = | French sentence 2 |
| English sentence 3 | = | |

# Visualization

*In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.*

The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

https://talktotransformer.com/

# Can BERT speak?

- Unified Language Model Pre-training for Natural Language Understanding and Generation
  - https://arxiv.org/abs/1905.03197
- BERT has a Mouth, and It Must Speak: BERT as a Markov Random Field Language Model
  - https://arxiv.org/abs/1902.04094
- Insertion Transformer: Flexible Sequence Generation via Insertion Operations
  - https://arxiv.org/abs/1902.03249
- Insertion-based Decoding with automatically Inferred Generation Order
  - https://arxiv.org/abs/1902.01370