# 人工智能技术及应用
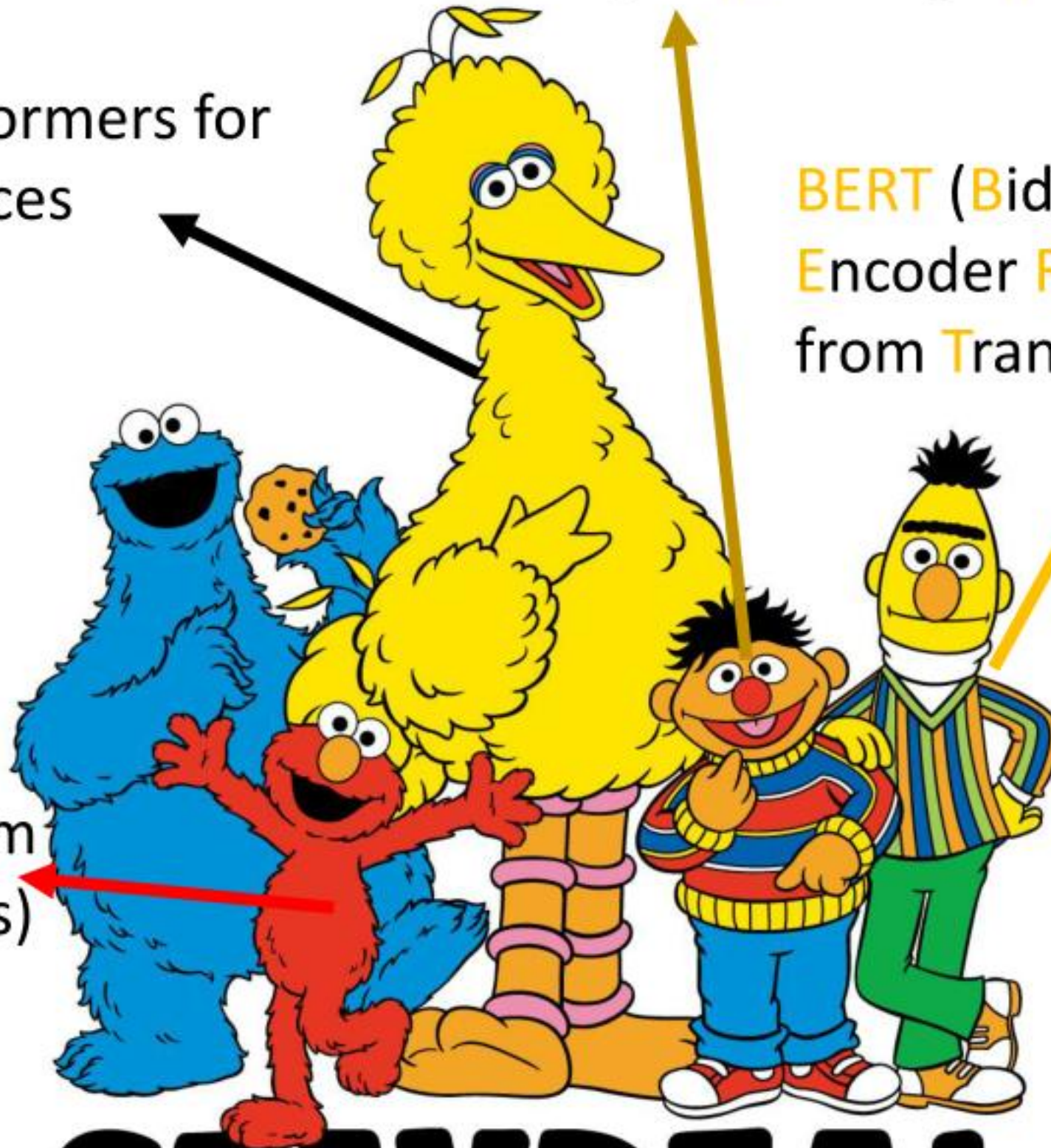## Artificial Intelligence and Application

Self-Supervised Learning

ERNIE (Enhanced Representation through Knowledge Integration)

Big Bird: Transformers for Longer Sequences

BERT (Bidirectional Encoder Representations from Transformers)

ELMo (Embeddings from Language Models)

STAYREAL

BERT

340M
parameters

Bertolt
Hoover

Source of image:
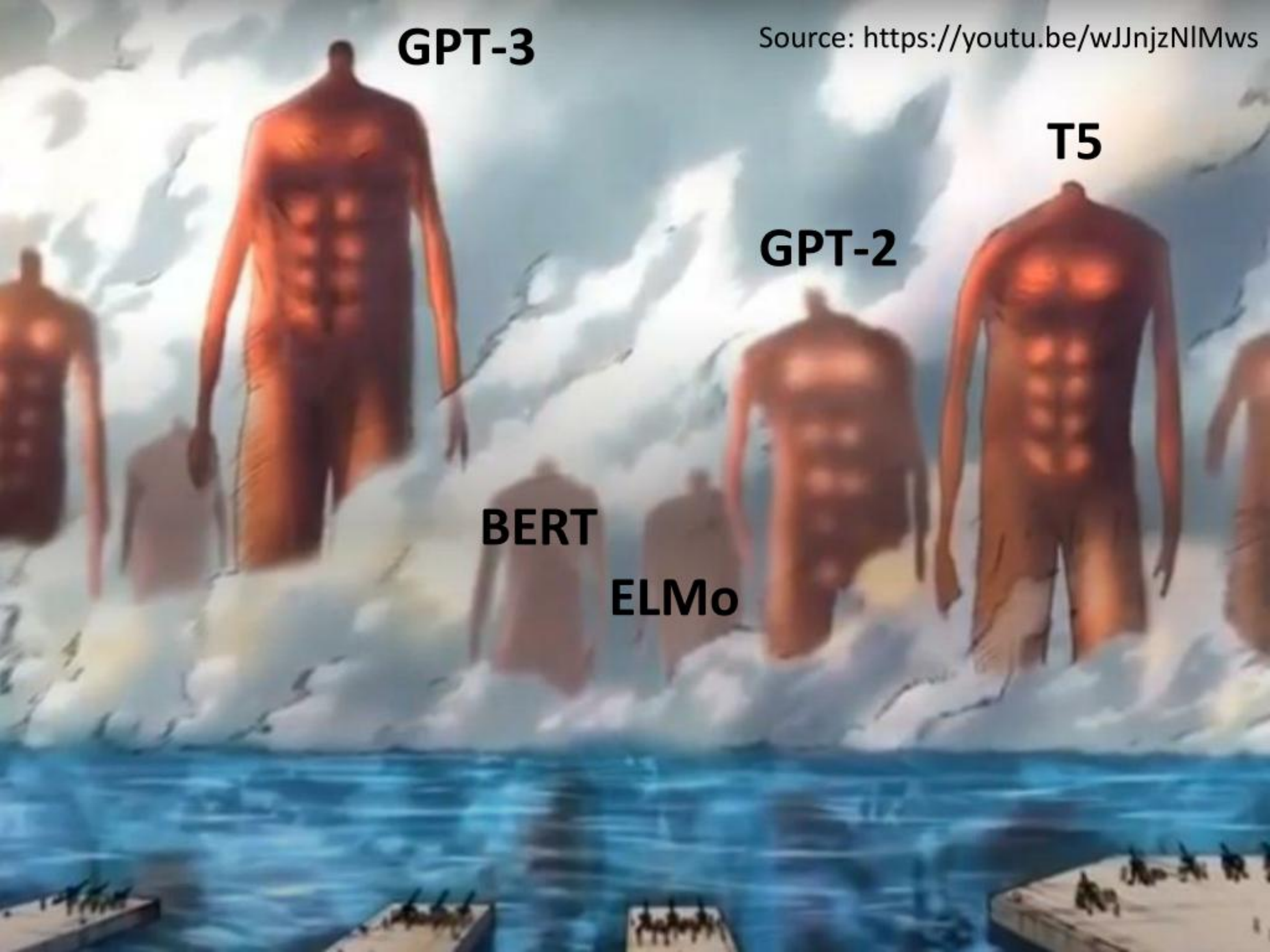https://leemeng.tw/attack_on_bert_transfer_learning_in_nlp.html

The models become larger
and larger ...

BERT
(340M)

ELMO
(94M)

GPT-2
(1542M)

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

The models become larger and larger ...

Turing NLG (17B)

GPT-3 is **10** times larger than Turing NLG. (175B)

GPT-2

Megatron (8B)

T5 (11B)

# Outline



BERT series
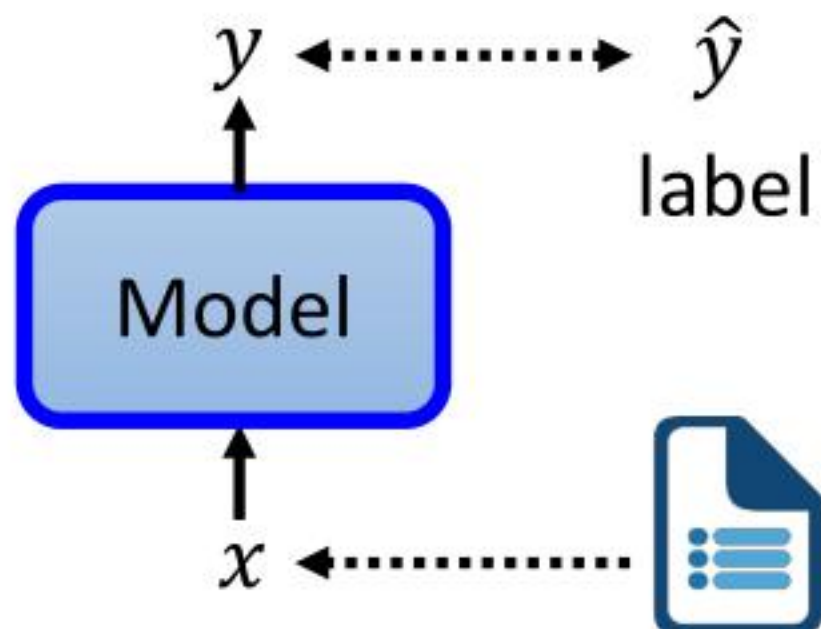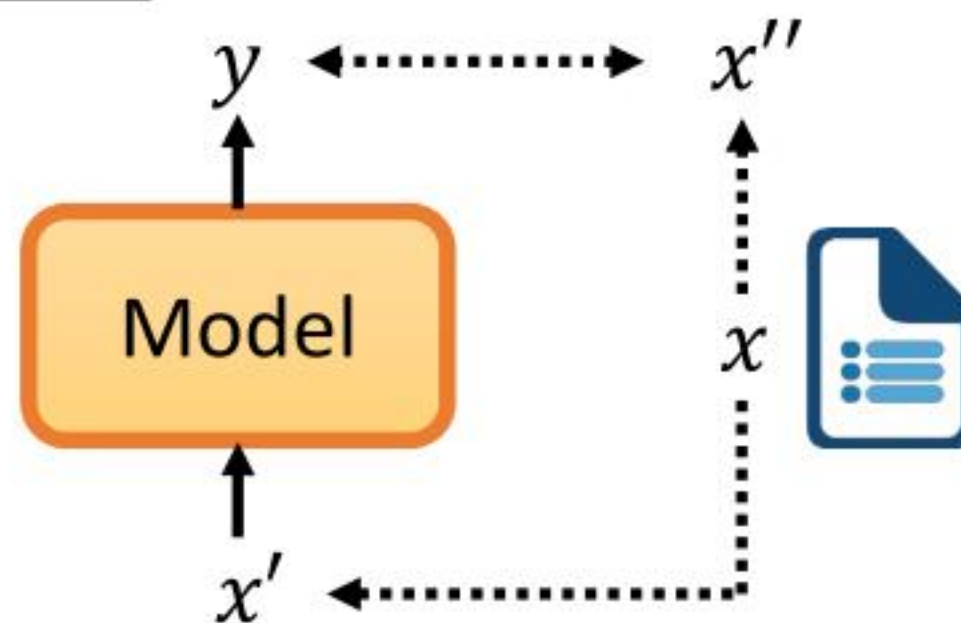
GPT series

# Self-supervised Learning

## Supervised

$$y \longleftrightarrow \hat{y}$$

label



Model

$x$

## Self-supervised

$$y \longleftrightarrow x''$$



Model

$x'$

$x$

---

**Yann LeCun**
2019年4月30日 · 🌐

<u>I now call it "self-supervised learning", because "unsupervised" is both a loaded and confusing term.</u>

In self-supervised learning, the system learns to predict part of its input from other parts of it input. In other words a portion of the input is used as a supervisory signal to a predictor fed with the remaining portion of the input.
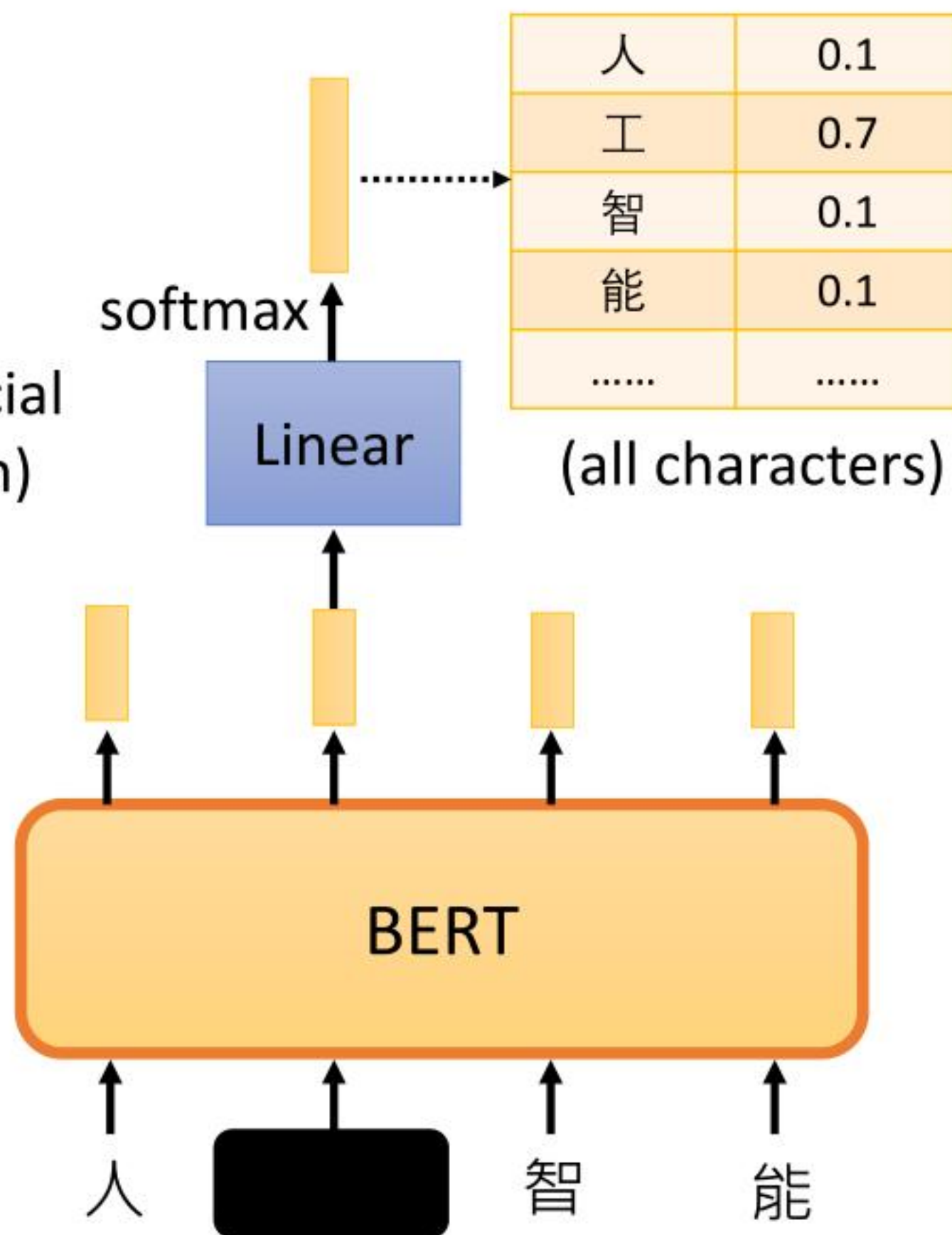
# Masking Input

# Masking Input

■ = MASK (special token)

or

■ = Random

一、天、大、小 ...

Transformer Encoder

Randomly masking some tokens

softmax

minimize cross entropy

Ground truth

工

Linear

BERT

人 ■ 智 能

# Next Sentence Prediction

Linear → Yes/No

BERT

[CLS]  $w_1$  $w_2$  [SEP]  $w_3$  $w_4$  $w_5$

Sentence 1          Sentence 2

- This approach is not helpful.

  Robustly optimized BERT approach
  (RoBERTa)  https://arxiv.org/abs/1907.11692

- **SOP**: Sentence order prediction

  Used in ALBERT
  https://arxiv.org/abs/1909.11942

Self-supervised Learning

**Pre-train**

BERT

- Masked token prediction
- Next sentence prediction

**Fine-tune**

Model for Task 1

Model for Task 2

Model for Task 3

*Downstream Tasks*

- The tasks we care
- We have a little bit labeled data.

# GLUE

General Language Understanding Evaluation (GLUE)

https://gluebenchmark.com/

- Corpus of Linguistic Acceptability (CoLA)
- Stanford Sentiment Treebank (SST-2)
- Microsoft Research Paraphrase Corpus (MRPC)
- Quora Question Pairs (QQP)
- Semantic Textual Similarity Benchmark (STS-B)
- Multi-Genre Natural Language Inference (MNLI)
- Question-answering NLI (QNLI)
- Recognizing Textual Entailment (RTE)
- Winograd NLI (WNLI)

GLUE also has Chinese version (https://www.cluebenchmarks.com/)

# BERT and its Family

- GLUE scores



Source of image: https://arxiv.org/abs/1905.00537

# How to use BERT – Case 1



class

Linear

Random initialization

**Better than random**
**Init by pre-train**

BERT

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

sentence

Input: sequence
output: class

Example:
Sentiment analysis

this is good
↓
positive

This is the model
to be learned.

# Pre-train v.s. Random Initialization

(fine-tune)                    (scratch)



Source of image: https://arxiv.org/abs/1908.05620

# How to use BERT – Case 2

class     class     class

Linear   Linear   Linear

BERT

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

sentence

Input: sequence
output: same as input

Example:
POS tagging

I  saw  a  saw

N  V  DET  N

# How to use BERT – Case 3

Input: two sequences
Output: a class

Example:
Natural Language Inferencee (NLI)

contradiction
entailment
neutral

Model

premise: A person on a horse
jumps over a broken down airplane

hypothesis: A person is at a diner.        contradiction

# How to use BERT – Case 3

# How to use BERT – Case 4

- Extraction-based Question Answering (QA)

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

**Query**: $Q = \{q_1, q_2, \cdots, q_M\}$
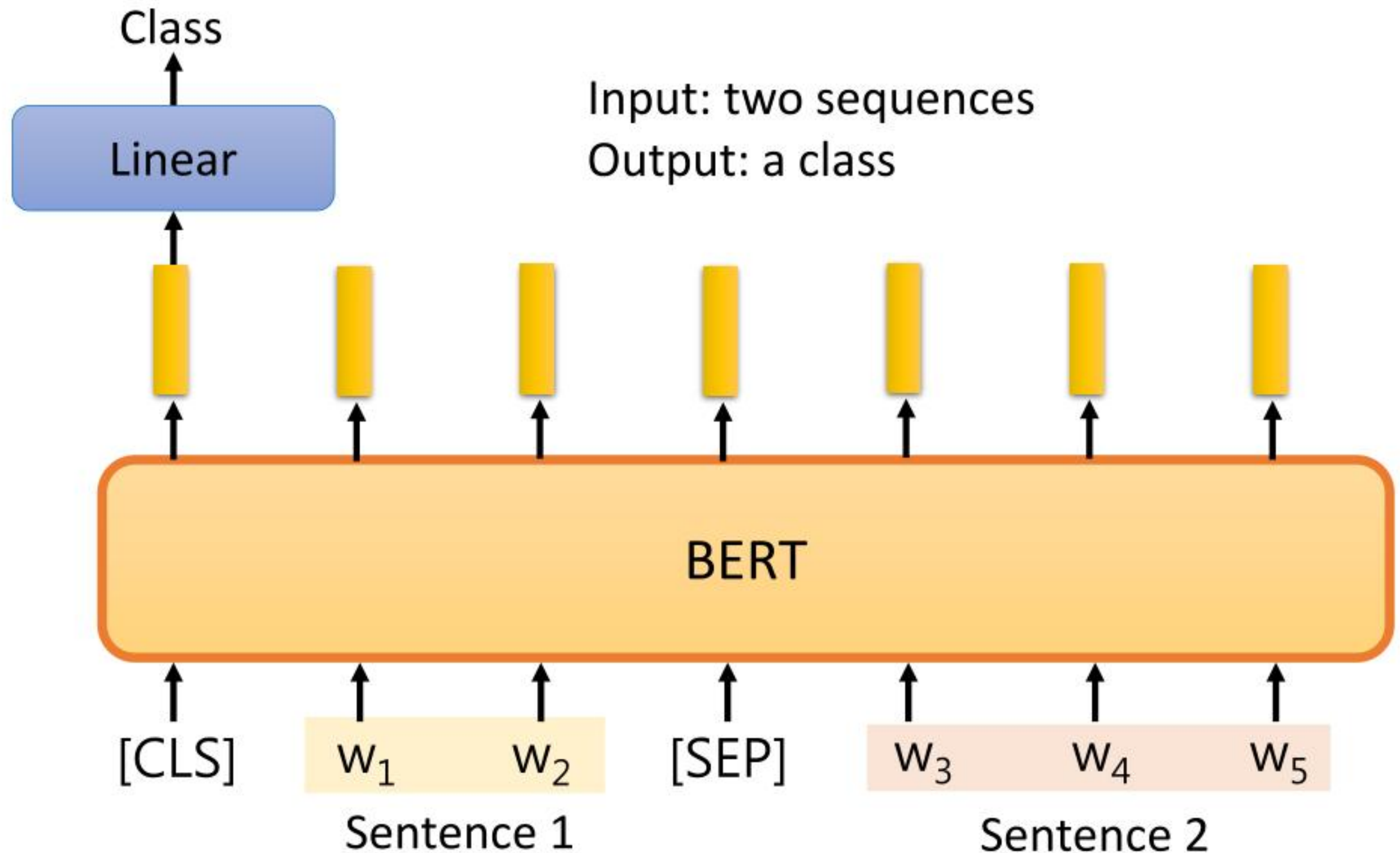
$$D \rightarrow \boxed{\text{QA Model}} \rightarrow s$$
$$Q \rightarrow \boxed{\text{QA Model}} \rightarrow e$$

output: two integers $(s, e)$

**Answer**: $A = \{d_s, \cdots, d_e\}$

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?
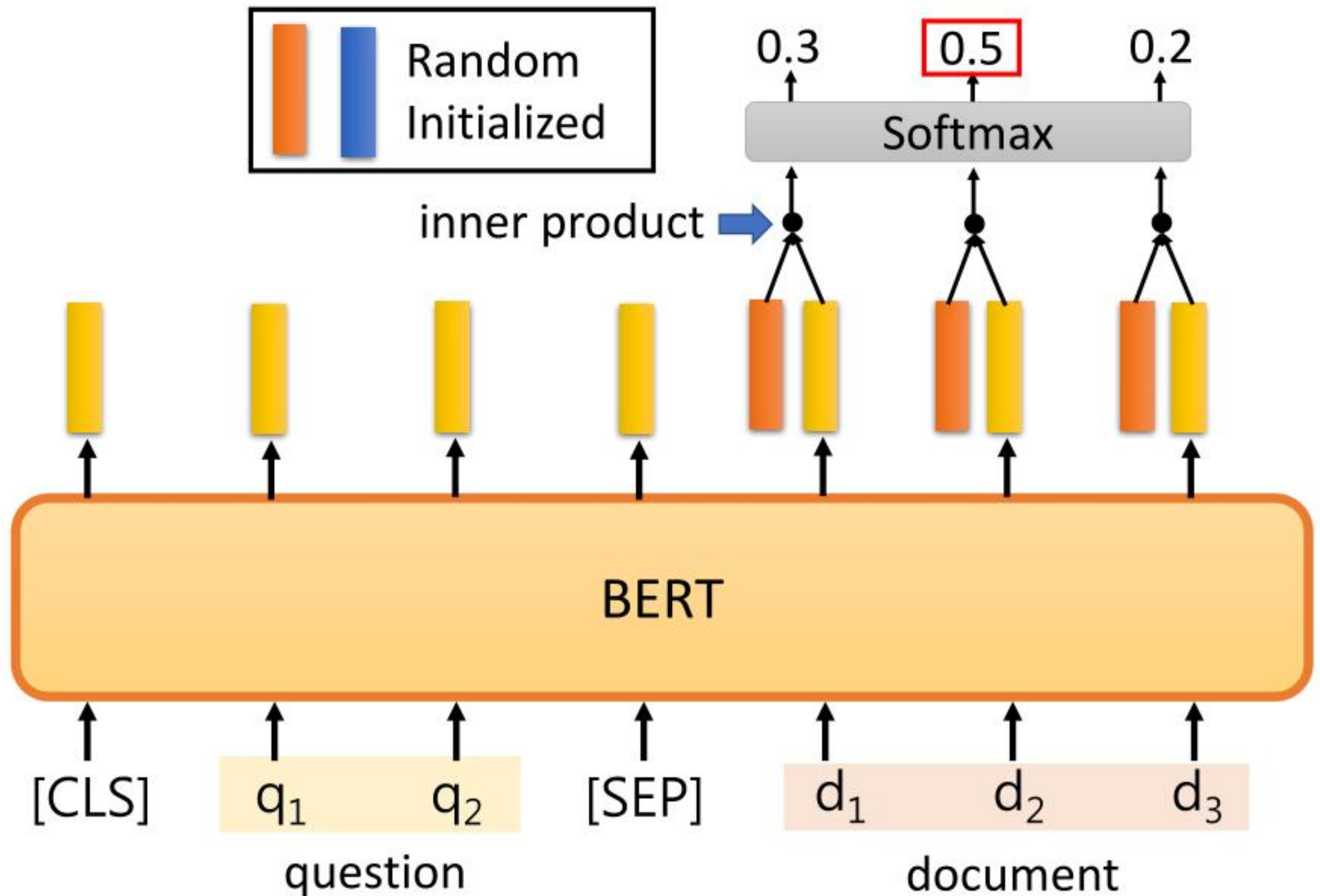**graupel**

Where do water droplets collide with ice crystals to form precipitation?
**within a cloud**      $s = 77, e = 79$

# How to use BERT – Case 4

# How to use BERT – Case 4

$s = 2$   $e = 3$

The answer is "$d_2 d_3$".

Random Initialized

inner product

0.1   0.2   0.7

Softmax

BERT

[CLS]   $q_1$   $q_2$   [SEP]   $d_1$   $d_2$   $d_3$
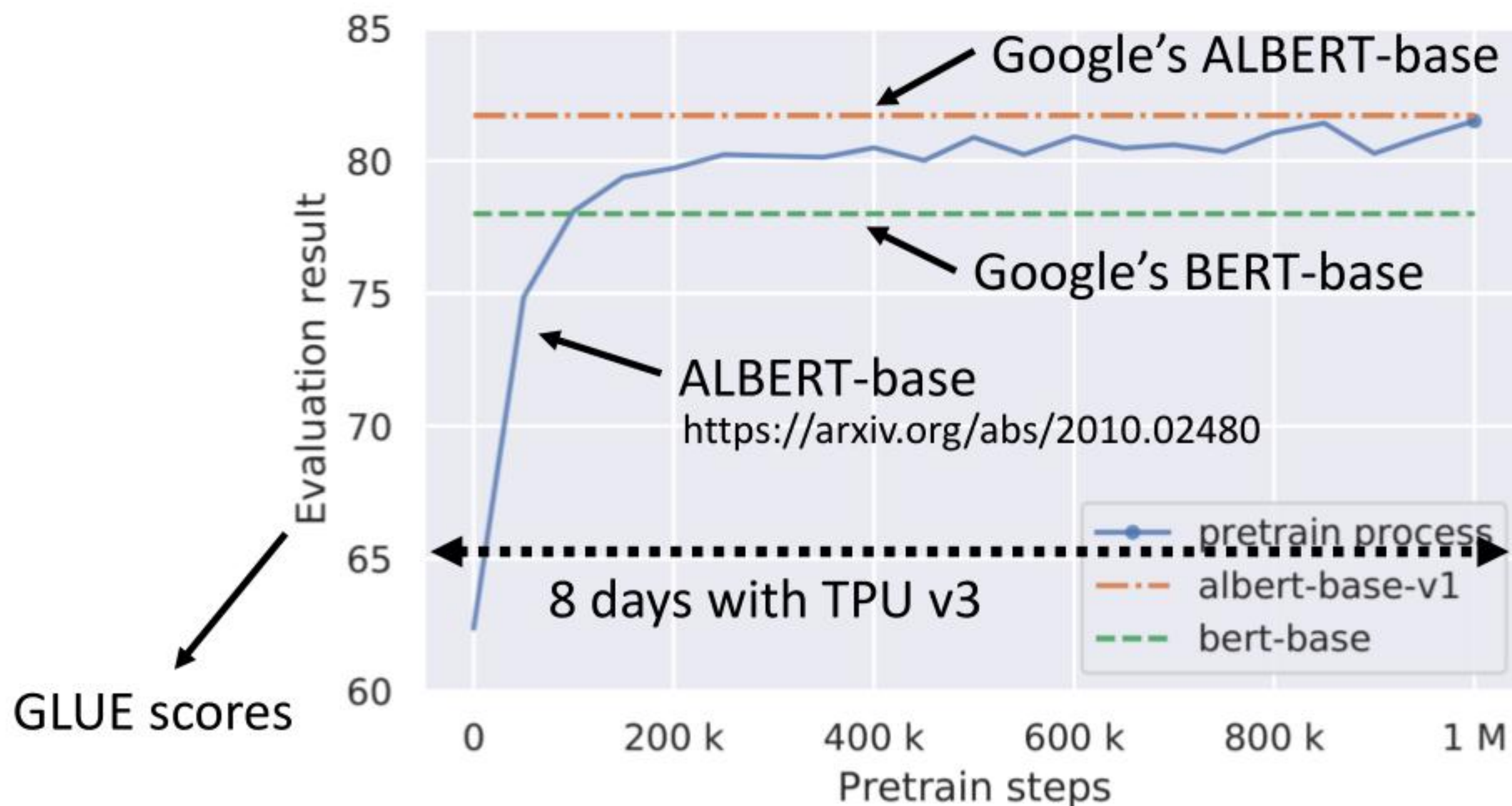
question

document

That's all!

# Training BERT is challenging!

Training data has more than **3 billions** of words.

**3000** times of **Harry Potter series**

# BERT Embryology (胚胎学)
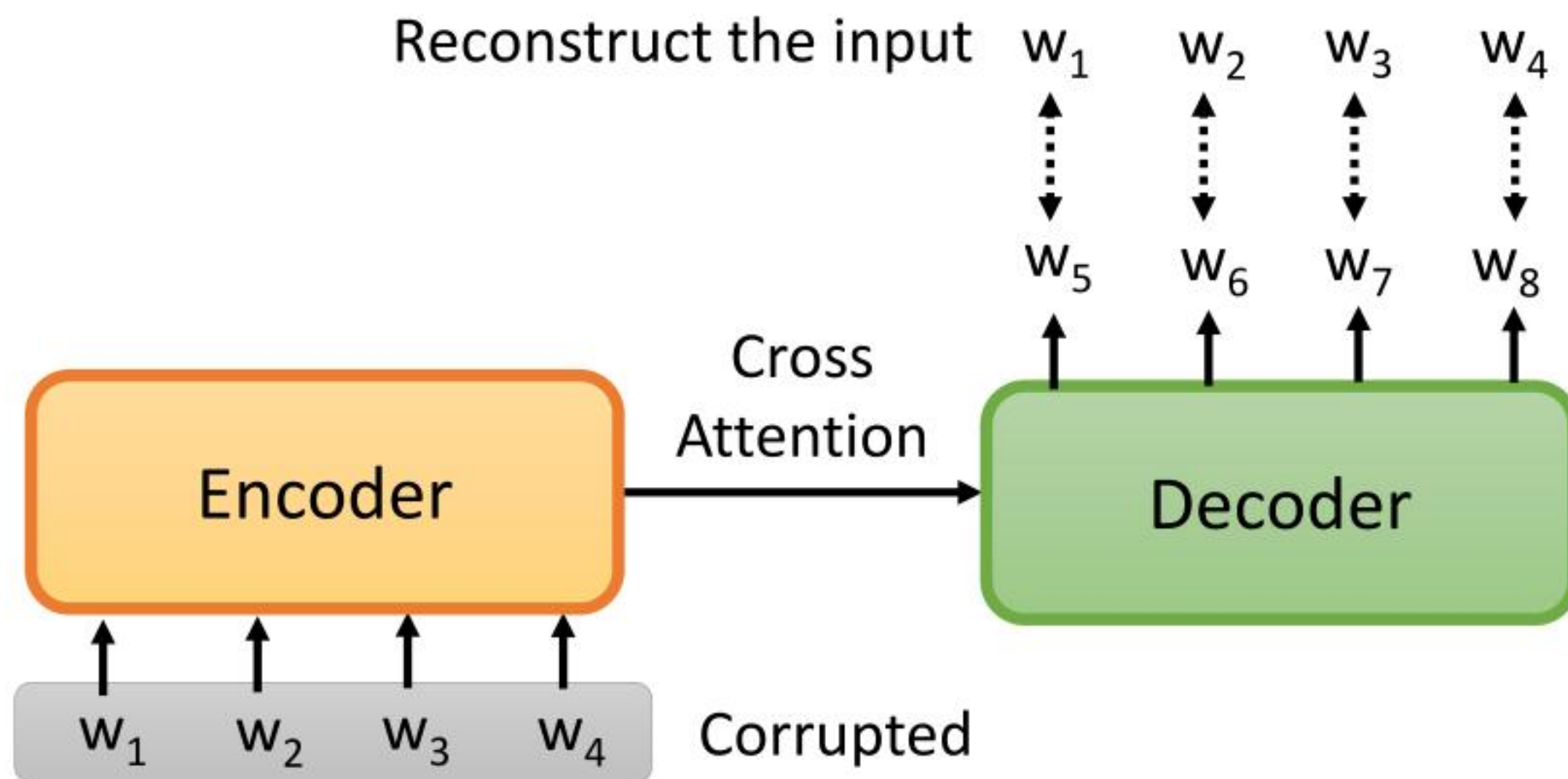
When does BERT know POS tagging,
syntactic parsing, semantics?

The answer is counterintuitive!

# Pre-training a seq2seq model

Reconstruct the input  $w_1$  $w_2$  $w_3$  $w_4$

$w_5$  $w_6$  $w_7$  $w_8$

Cross
Attention

Encoder

Decoder

$w_1$  $w_2$  $w_3$  $w_4$  Corrupted

# MASS / BART

MASS

BART

A B [SEP] C D E

A B [SEP] C ▢ E

A B [SEP] C E
(Delete "D")

C D E [SEP] A B
(permutation)

D E A B [SEP] C
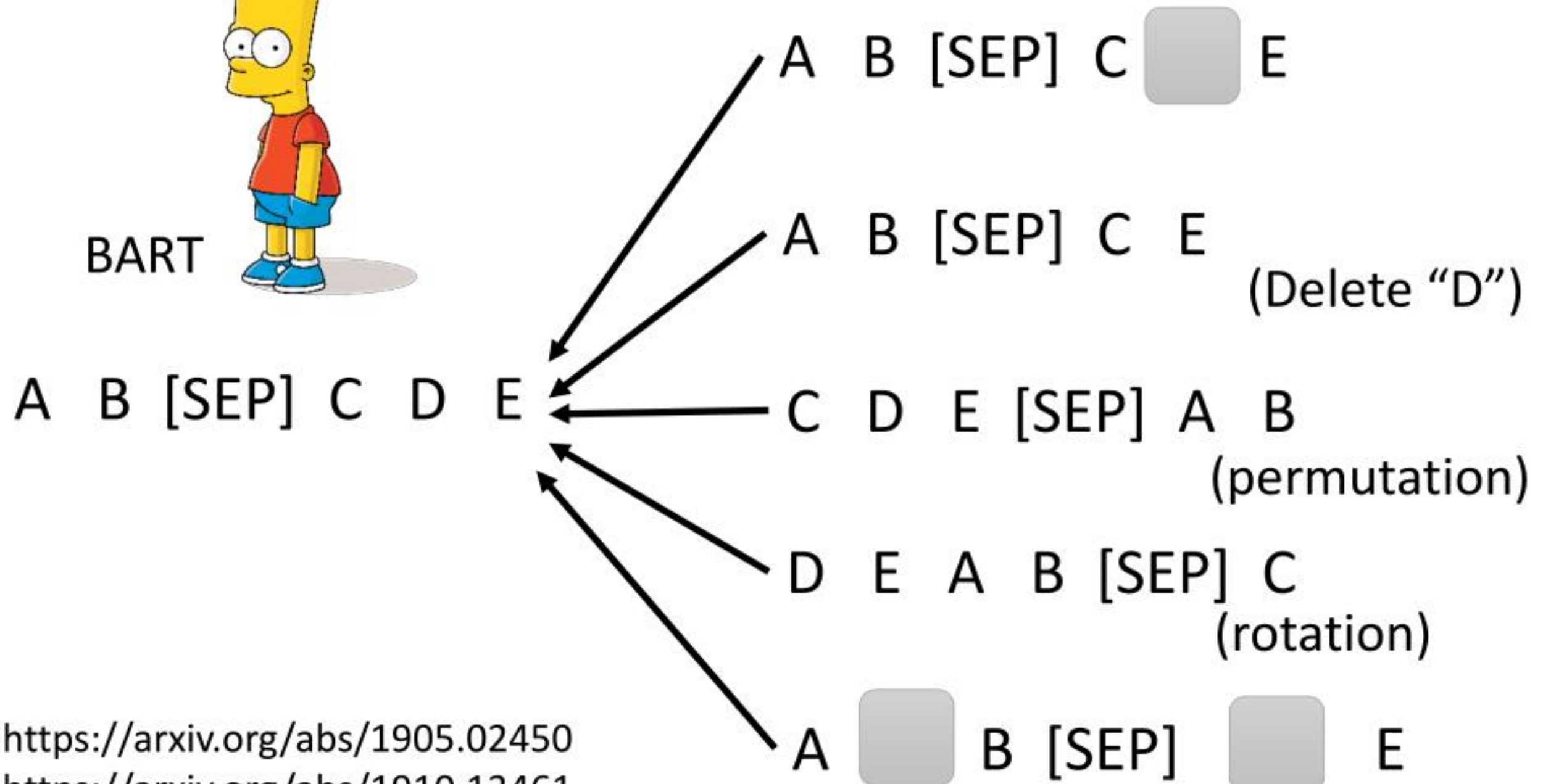(rotation)

A ▢ B [SEP] ▢ E

**Text Infilling**

# T5 – Comparison

- Transfer Text-to-Text Transformer (T5)
- Colossal Clean Crawled Corpus (C4)

| Objective | Inputs | Targets |
|---|---|---|
| Prefix language modeling | Thank you for inviting | me to your party last week . |
| BERT-style | Thank you \<M> \<M> me to your party apple week . | (original text) |
| Deshuffling | party me for your to . las | |
| I.i.d. noise, mask tokens | Thank you \<M> \<M> me t | |
| I.i.d. noise, replace spans | Thank you \<X> me to yo | |
| I.i.d. noise, drop tokens | Thank you me to your pa | |
| Random spans | Thank you \<X> to \<Y> we | |

# Why does BERT work?

embedding

Represent the **meaning** of "智"

The tokens with similar meaning have similar embedding.



BERT

人　工　智　能

吃苹果　　草

电　鸟　鱼　苹果手机

**Context is considered.**

# Why does BERT work?

compute cosine similarity

BERT

BERT

喝　苹　果　汁　　苹　果　电　脑

Cosine Similarities of BERT Embeddings

今天买了苹果来吃

进口苹果平均每公斤下跌12.3%

苹果茶真难喝
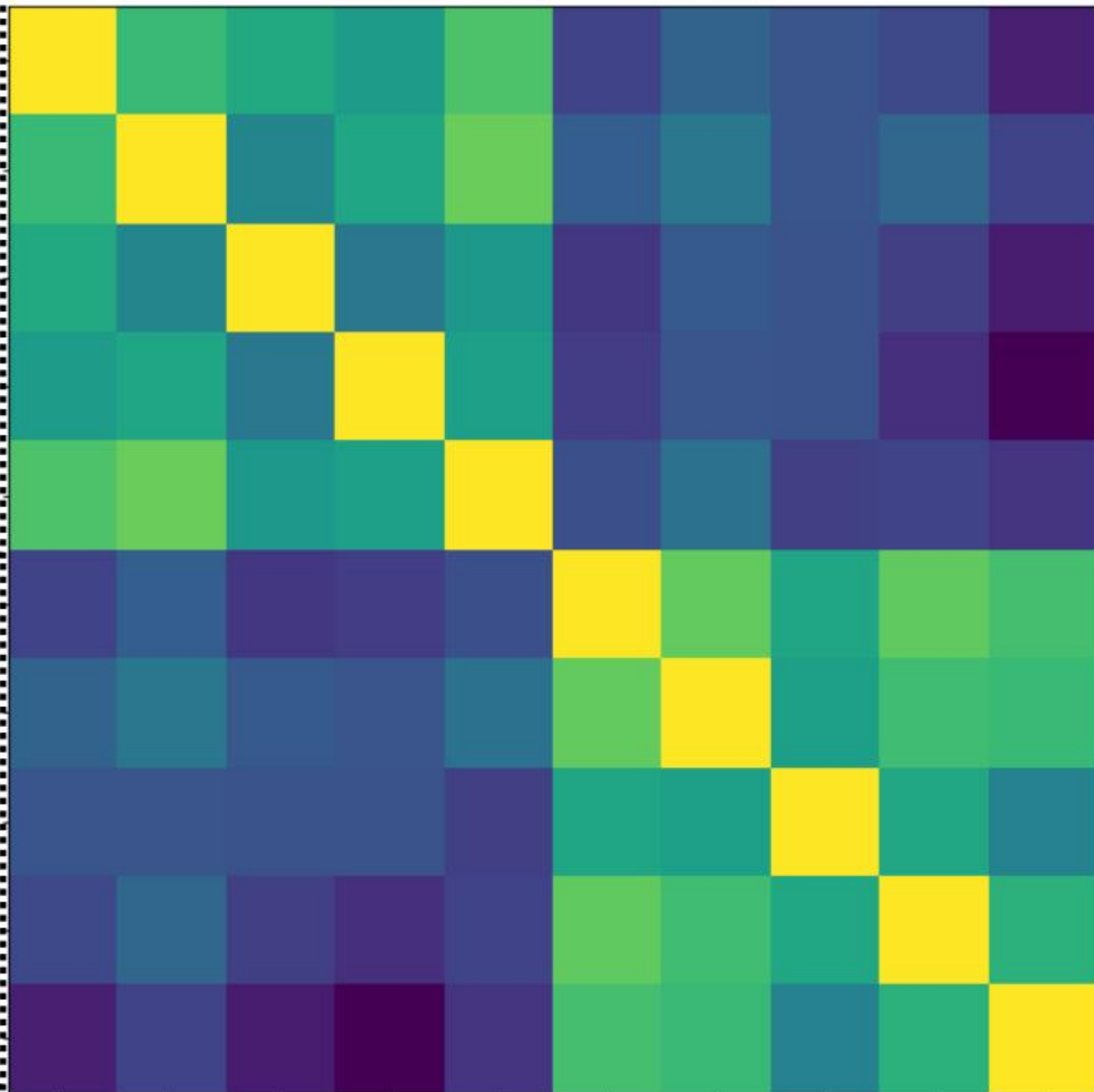
苹果收获的季节快到了

进口苹果因防止水分流失添加人工果蜡

苹果即将于下月发布新款iPhone

苹果获新Face ID专利

今天买了苹果手机

苹果的股价又跌了

苹果指纹识别技术

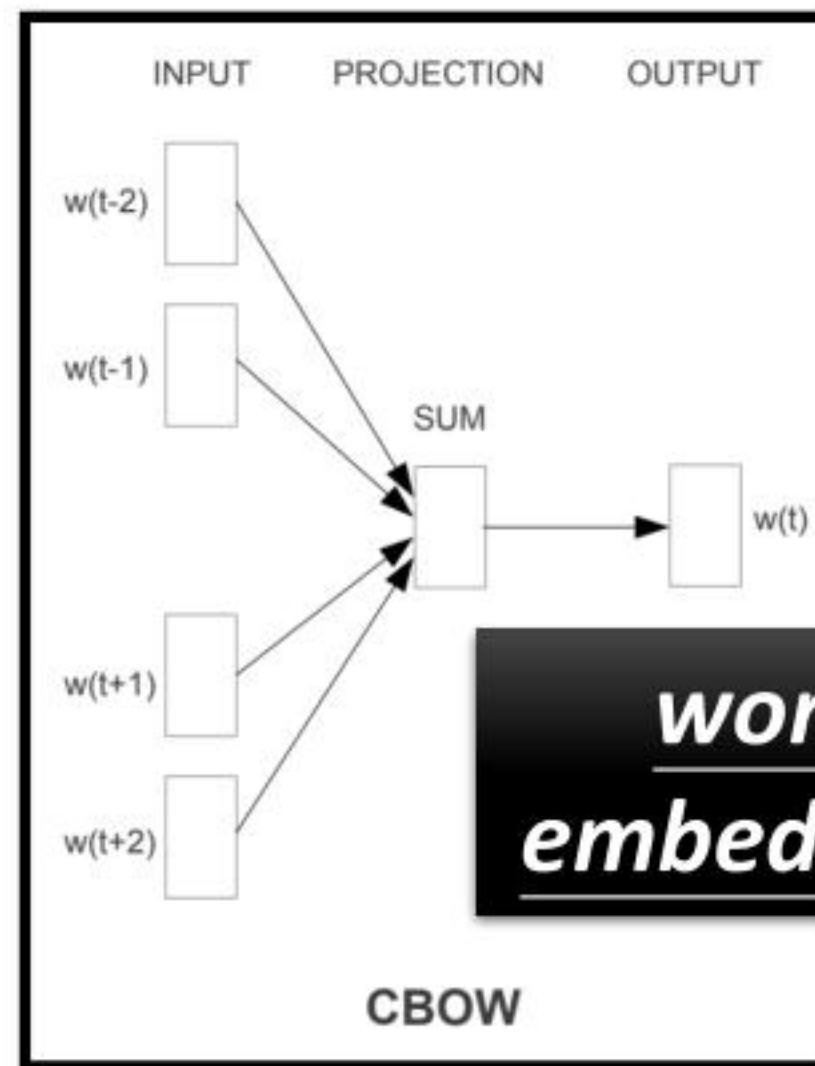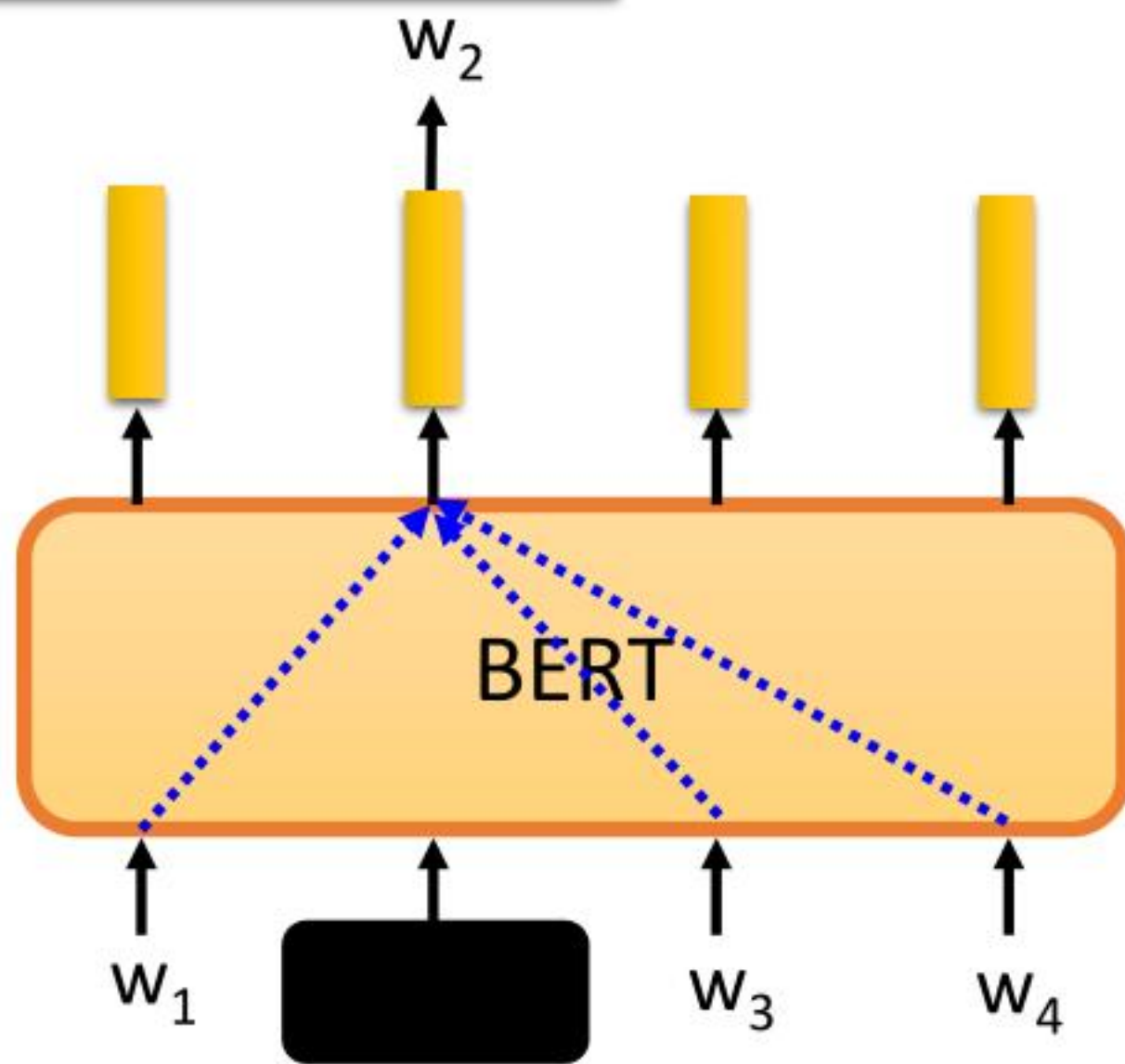# Why does BERT work?

**Contextualized word embedding**

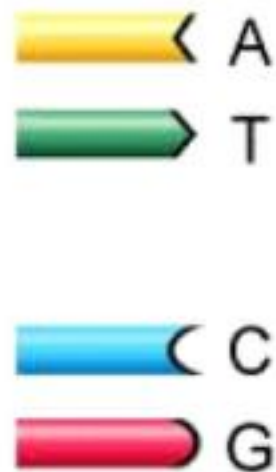You shall know a word by the company it keeps

John Rupert Firth

$w_2$

BERT

$w_1$ $w_3$ $w_4$

INPUT     PROJECTION     OUTPUT

w(t-2)

w(t-1)

SUM

w(t)

w(t+1)

w(t+2)

**word embedding**

CBOW

# Why does BERT work?

- Applying BERT to **protein**, **DNA**, **music classification**

| class | DNA sequence |
|---|---|
| EI | CCAGCTGCATCACAGGAGGCCAGCG |
| EI | AGACCCGCCGGGAGGCGGAGGACC |
| IE | AACGTGGCCTCCTTGTGCCCTTCCC |
| IE | CCACTCAGCCAGGCCCTTCTTCTCCT |
| IE | CCTGATCTGGGTCTCCCCTCCCACCC |
| IE | AGCCCTCAACCCTTCTGTCTCACCCT |
| IE | CCACTCAGCCAGGCCCTTCTTCTCCT |
| N | CTGTGTTCACCACATCAAGCGCCGGG |
| N | GTGTTACCGAGGGCATTTCTAACAGT |
| N | TCTGAGCTCTGCATTTGTCTATTCTCC |

# *Why does BERT work?*

class

Linear ┄┄▶ Random initialization

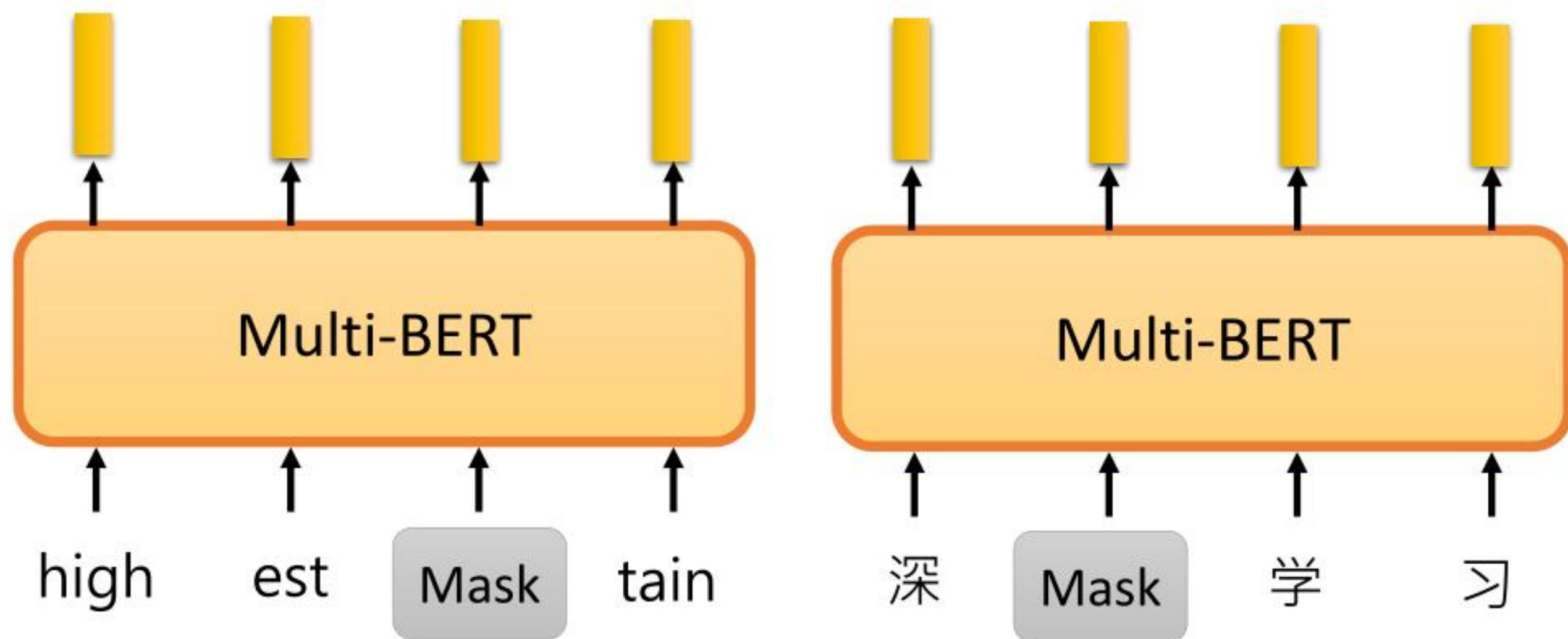pre-train on English

Init by pre-train

BERT

| A | we |
|---|-----|
| T | you |
| C | he |
| G | she |

[CLS]  we  she  we  he

DNA sequence ┄┄▶ A  G  A  C

# Why does BERT work?

- Applying BERT to **protein**, **DNA**, **music classification**

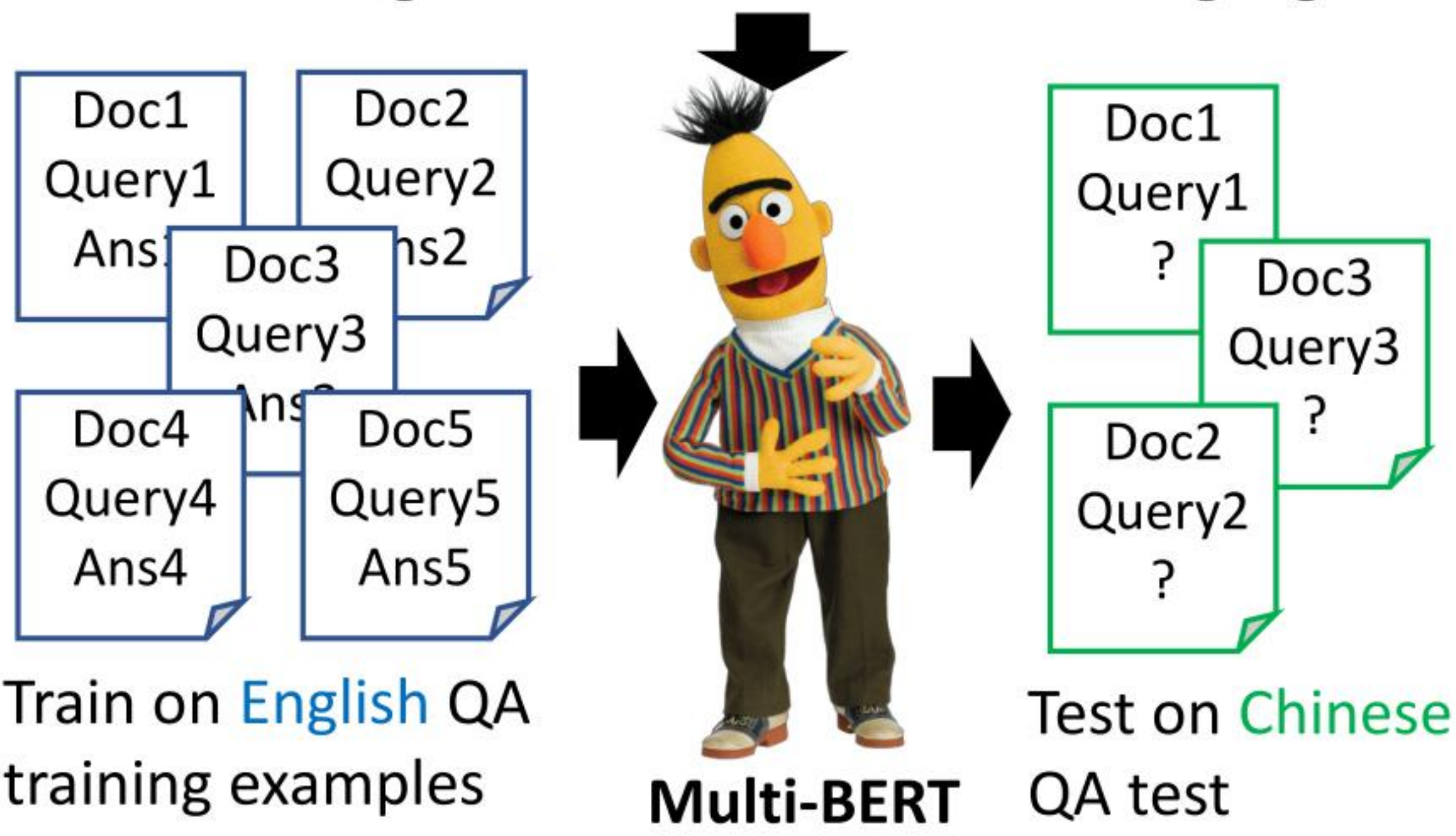| | Protein | | | DNA | | | | Music |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | localization | stability | fluorescence | H3 | H4 | H3K9ac | Splice | composer |
| specific | 69.0 | 76.0 | 63.0 | 87.3 | 87.3 | 79.1 | 94.1 | - |
| BERT | 64.8 | 74.5 | 63.7 | 83.0 | 86.2 | 78.3 | 97.5 | 55.2 |
| re-emb | 63.3 | 75.4 | 37.3 | 78.5 | 83.7 | 76.3 | 95.6 | 55.2 |
| rand | 58.6 | 65.8 | 27.5 | 75.6 | 66.5 | 72.8 | 95 | 36 |

# Multi-lingual BERT



Training a BERT model by many different languages.

# Zero-shot Reading Comprehension
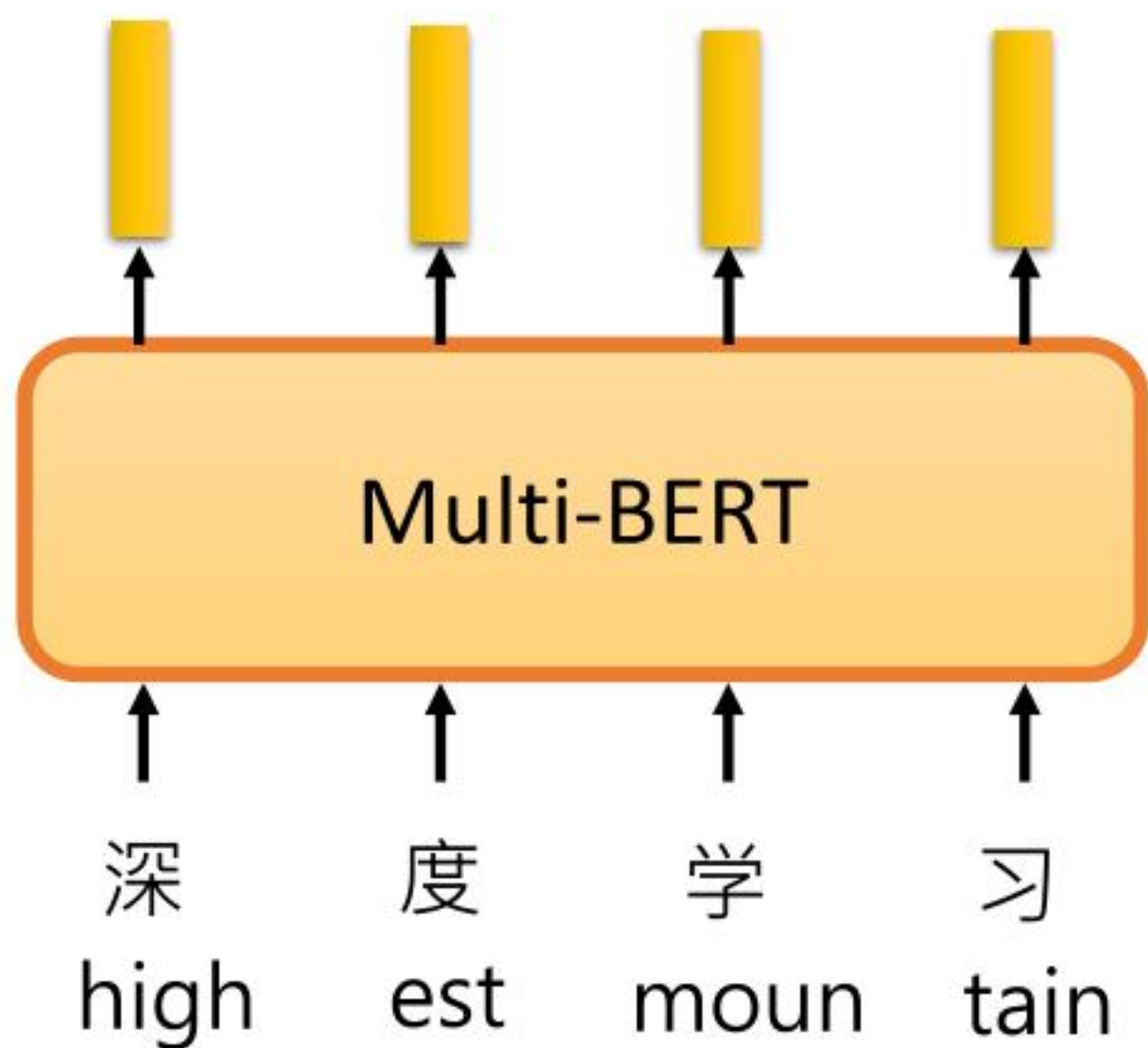
Training on the sentences of 104 languages

Doc1
Query1
Ans1

Doc2
Query2
Ans2

Doc3
Query3
Ans3

Doc4
Query4
Ans4

Doc5
Query5
Ans5

Train on English QA training examples

**Multi-BERT**

Doc1
Query1
?

Doc3
Query3
?

Doc2
Query2
?

Test on Chinese QA test

# Zero-shot Reading Comprehension

- English: SQuAD, Chinese: DRCD

| Model | Pre-train | Fine-tune | Test | EM | F1 |
|-------|-----------|-----------|------|-----|-----|
| QANet | none | Chinese | Chinese | 66.1 | 78.1 |
| BERT | Chinese | Chinese | | 82.0 | 89.1 |
| | 104 languages | Chinese | | 81.2 | 88.7 |
| | | English | | 63.3 | 78.8 |
| | | Chinese + English | | 82.6 | 90.1 |

F1 score of Human performance is 93.30%

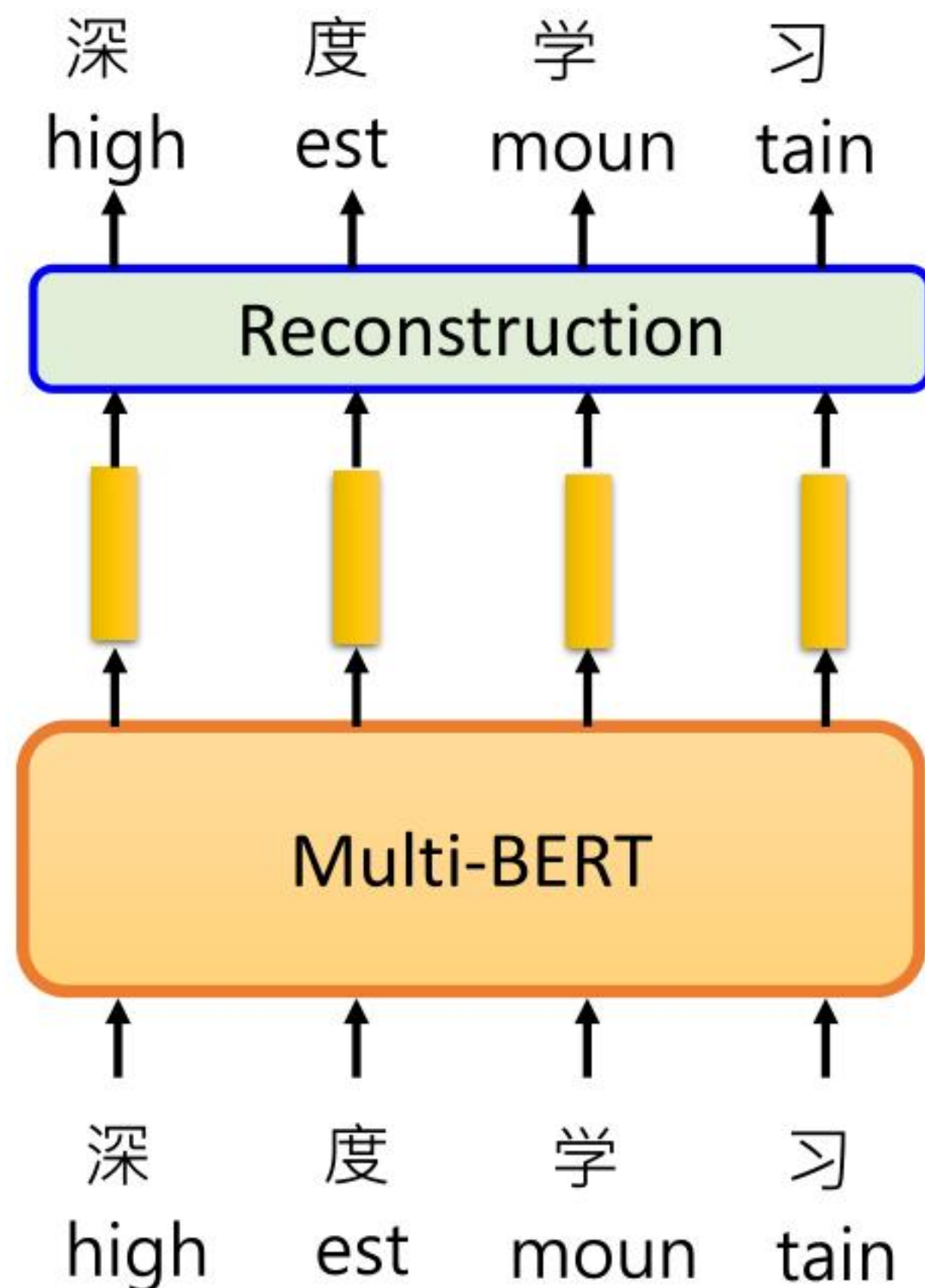# Cross-lingual Alignment?

## *Weird???*

跳 jump

游 swim

兔 rabbit

鱼 fish
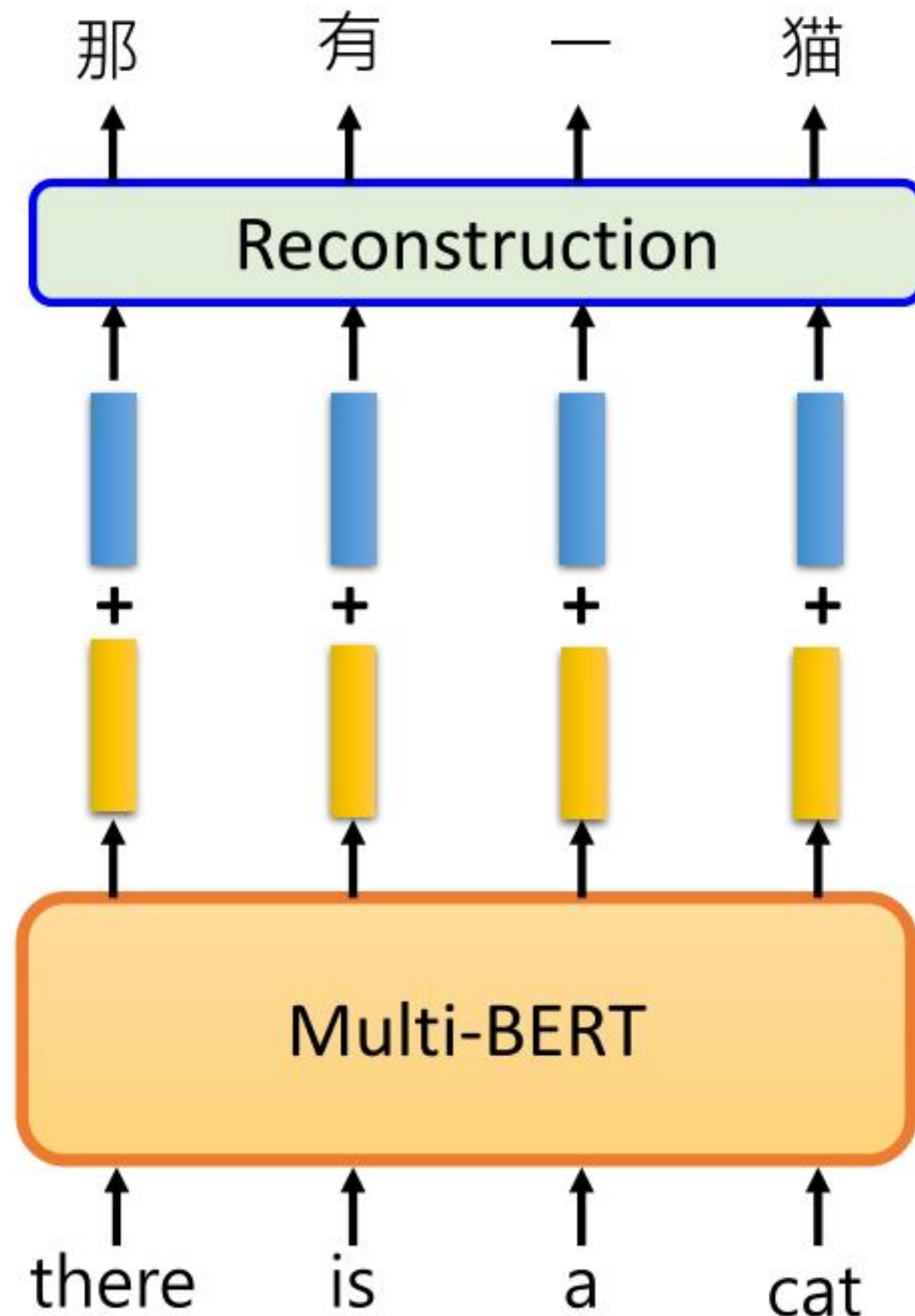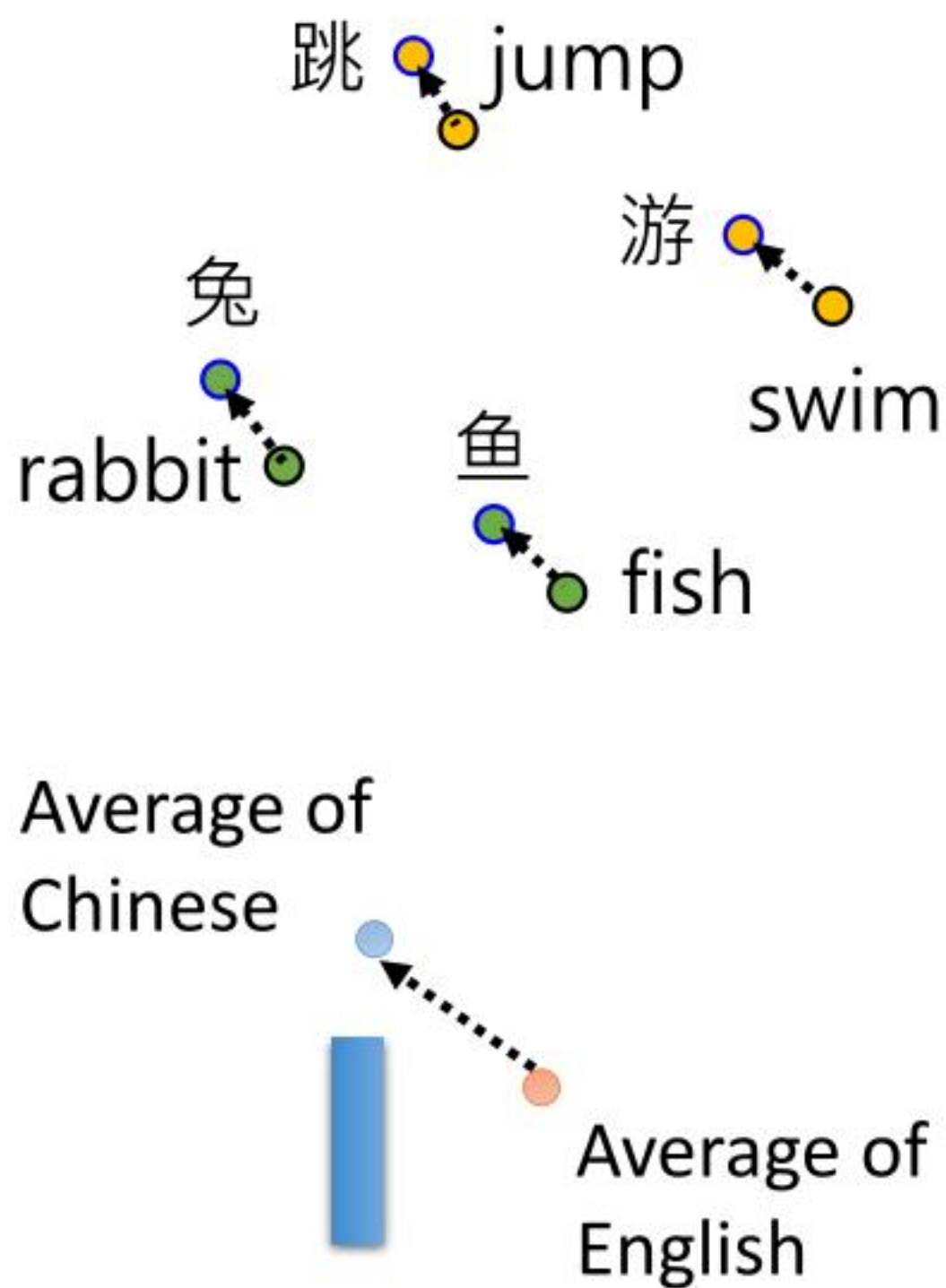
If the embedding is language independent ...

How to correctly reconstruct?

There must be language information.

深 度 学 习
high est moun tain

Reconstruction

Multi-BERT

深 度 学 习
high est moun tain
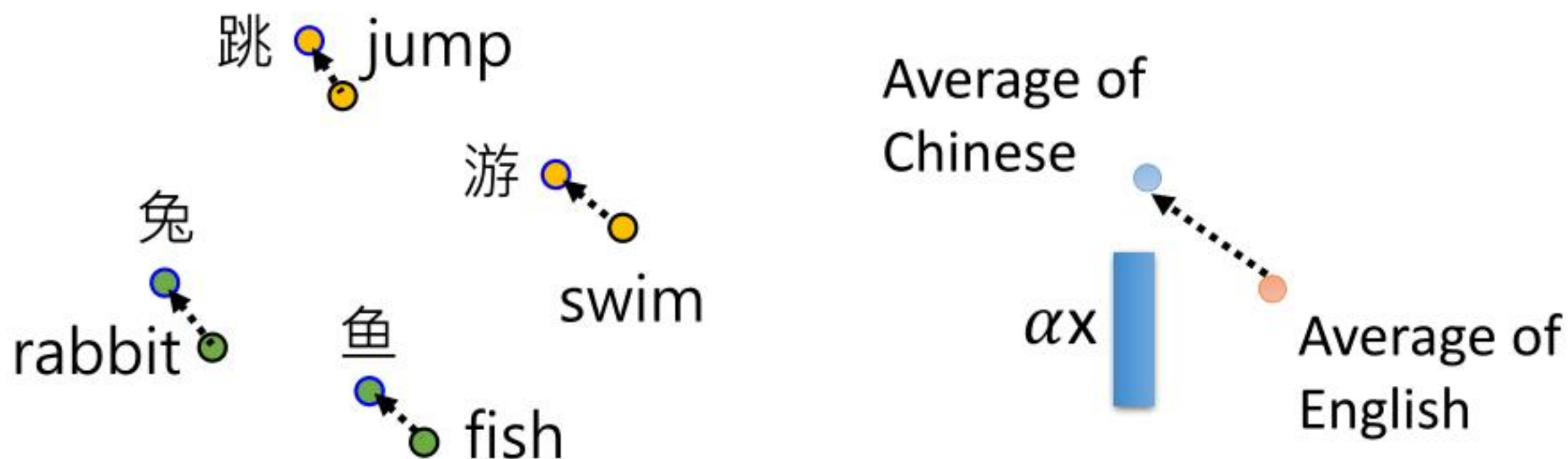
# If this is true ...

| Input (en) | The girl that can help me is all the way across town. There is no one who can help me. |
| --- | --- |
| Ground Truth (zh) | 能帮助我的女孩在小镇的另一边。没有人能帮助我。。 |
| en→zh, $\alpha = 1$ | . 孩，can 来我是all the way across 市。。 There 是无人人can help 我。 |
| en→zh, $\alpha = 2$ | . 孩的的家我是这个人的市。。他是他人人的到我。 |
| en→zh, $\alpha = 3$ | 。，的的的他是的个的的，。：他是他人，的。他。 |

Unsupervised token-level translation ☺