# Natural Language Processing

## Lecture 13: Contextualized Representation (Pre-trained Models)

Amirkabir University of Technology

Dr Momtazi

# Outline
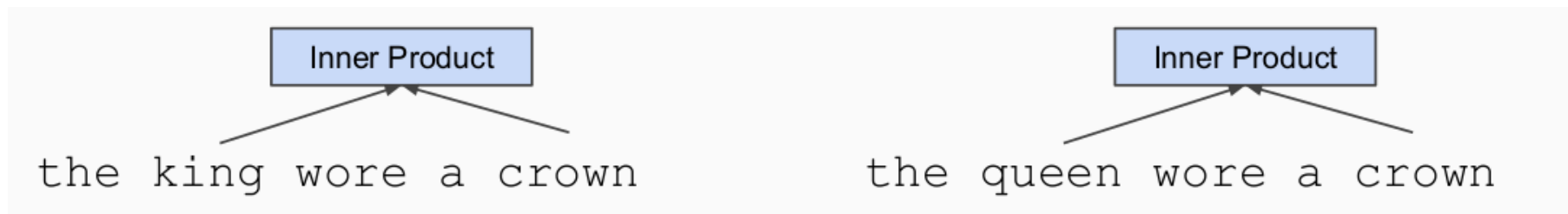
- **Introduction**
- ELMO
- BERT

# Pre-training in NLP

- Word embeddings are the basis of deep learning for NLP

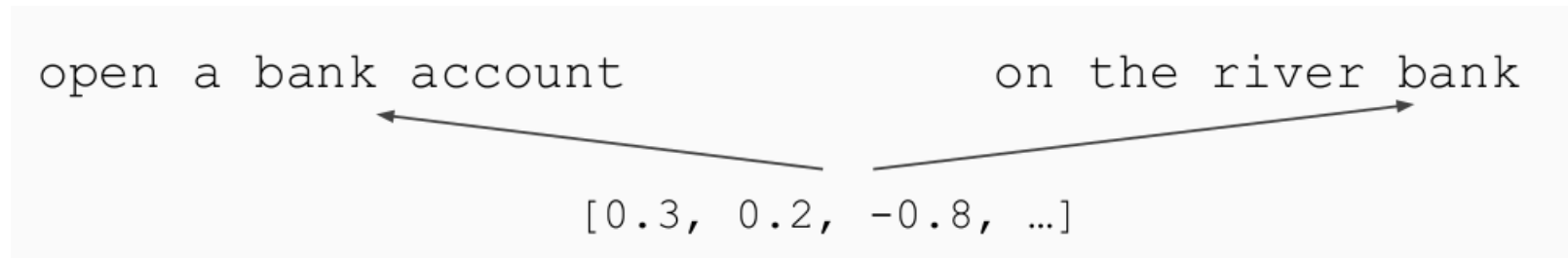| king | queen |
|------|-------|
| ↓ | ↓ |
| [-0.5, -0.9, 1.4, …] | [-0.6, -0.8, -0.2, …] |

- Word embeddings (word2vec, GloVe) are often *pre-trained* on text corpus from co-occurrence statistics

# Contextual Representations

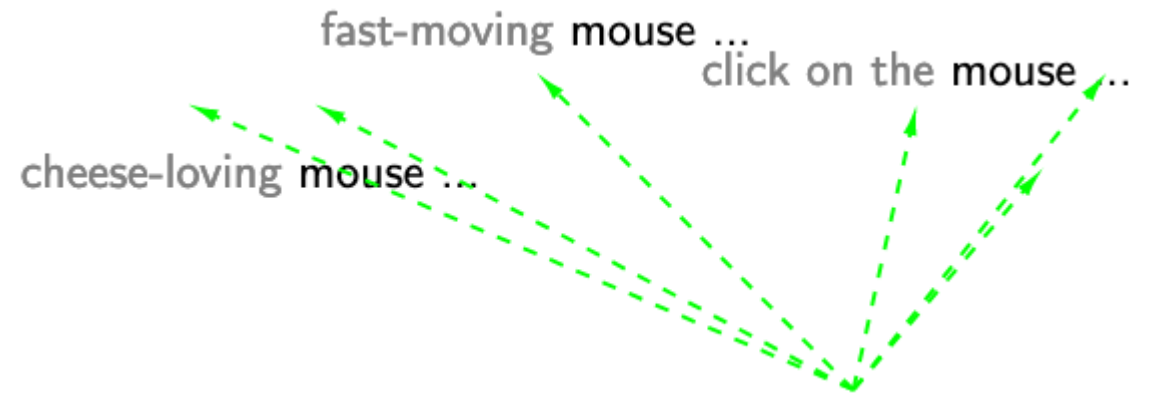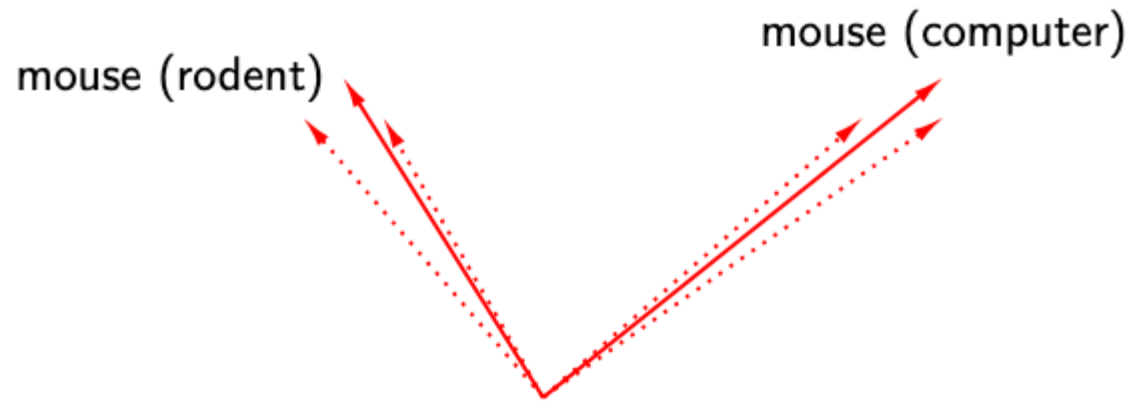- **Problem**: Word embeddings are applied in a context free manner

```
open a bank account                    on the river bank
                    [0.3, 0.2, -0.8, …]
```

- **Solution**: Train *contextual* representations on text corpus

```
  [0.9, -0.2, 1.6, …]                              [-1.9, -0.4, 0.1, …]
          ↑                                                  ↑
  open a bank account                              on the river bank
```

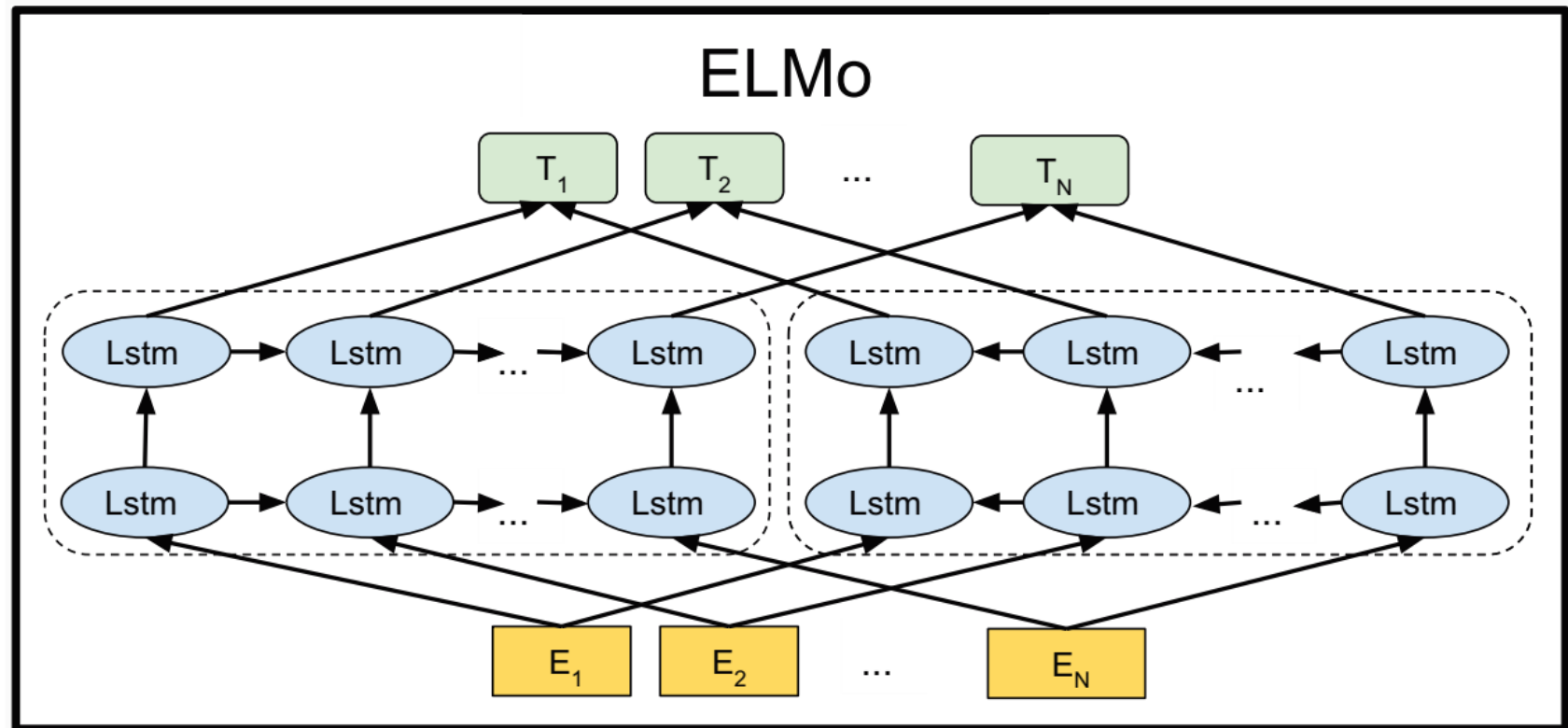# Static vs Contextualized Representation

# Outline

- Introduction
- **ELMO**
- BERT

# ELMO

- Bidirectional language model
- Train Separate Left-to-Right and Right-to-Left LMs

# ELMO

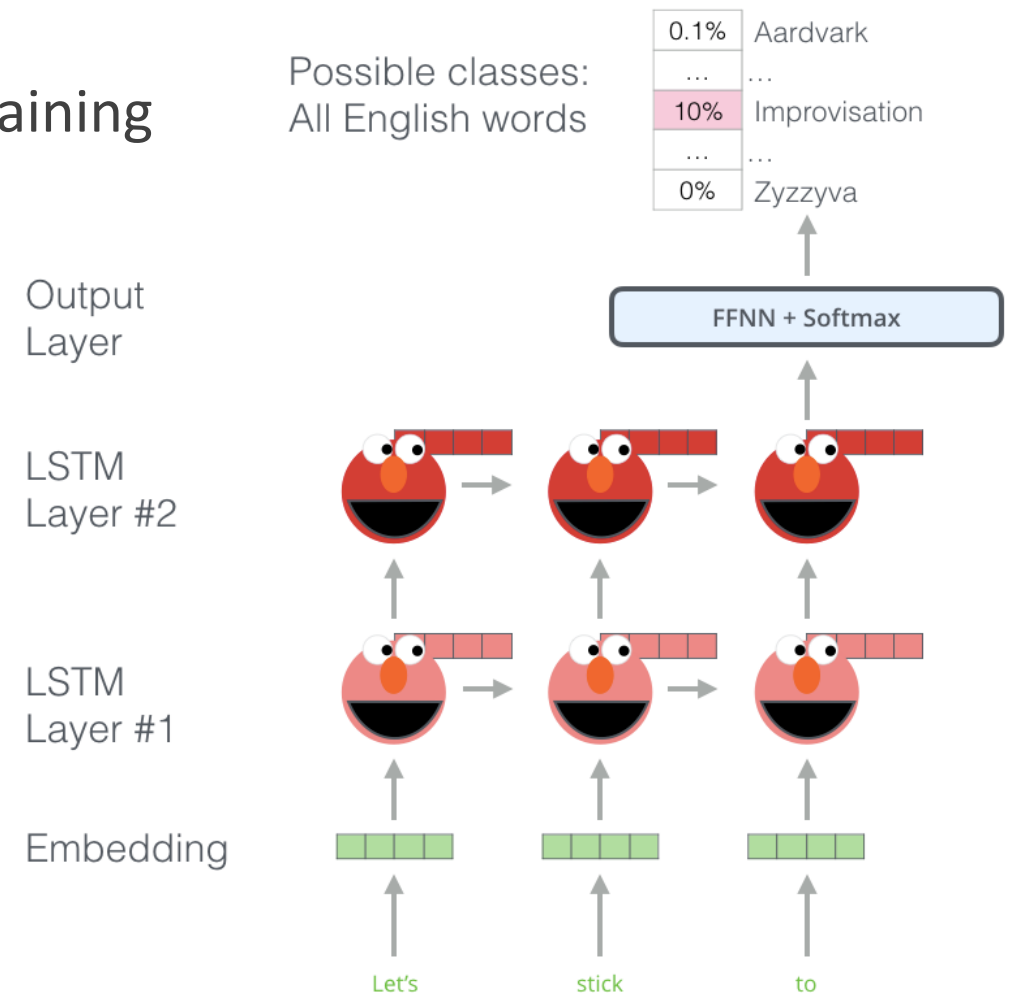- Bidirectional training

- Forward language model

$$p(t_1, t_2;, \dots, t_N) =| \prod_{k=1}^{N} p(t_k|t_1, t_2;, \dots, t_{k-1})$$

- Backward language model

$$p(t_1, t_2;, \dots, t_N) = \prod_{k=1}^{N} p(t_k|t_{k+1}, t_{k+2};, \dots, t_N)$$

# ELMO

- Predicting the next word in forward training

- Predicting the previous word in backward training

# ELMO

- Extracting embedding



1- Concatenate hidden layers
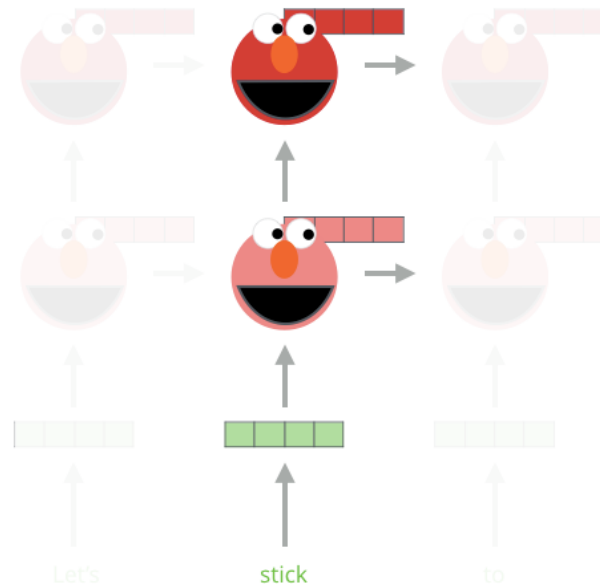
2- Multiply each vector by a weight based on the task

$x \quad s_2$

$x \quad s_1$
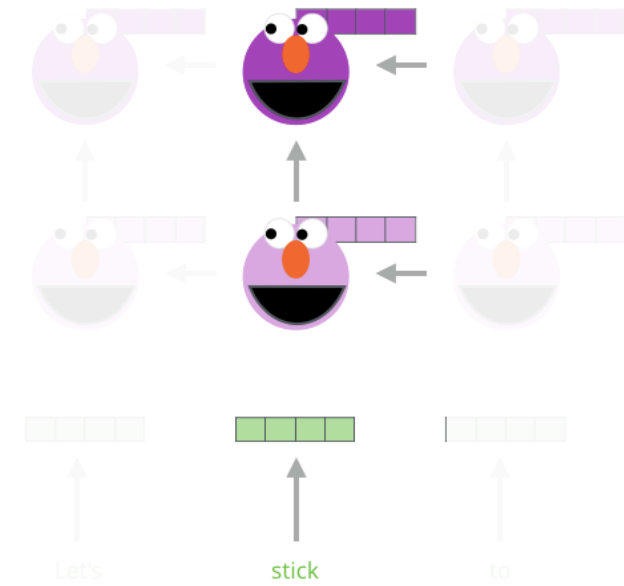
$x \quad s_0$

3- Sum the (now weighted) vectors

ELMo embedding of "stick" for this task in this context

Forward Language Model

Backward Language Model

Let's          stick          to

Let's          stick          to

# ELMO

- Extracting embedding

ELMo is a task specific representation. A down-stream task learns weighting parameters

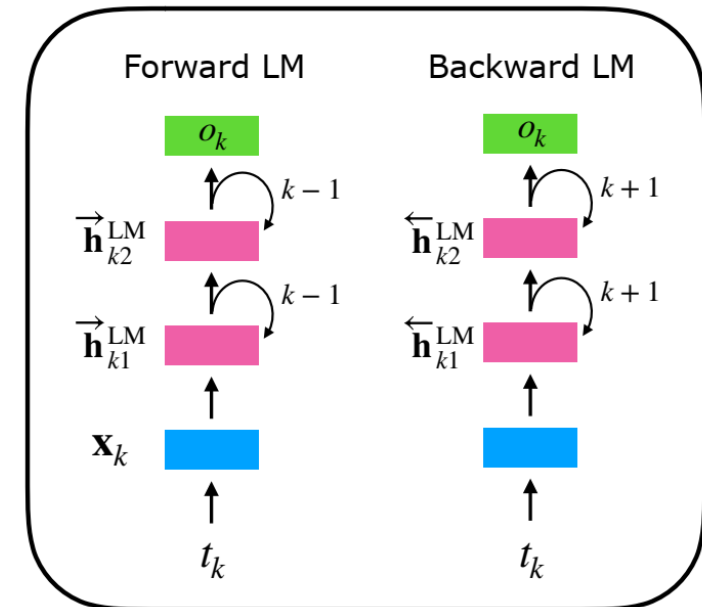$$\mathbf{ELMo}_k^{\text{task}} = \gamma^{\text{task}} \times \sum \begin{cases} s_2^{\text{task}} \times \mathbf{h}_{k2}^{\text{LM}} \\ \\ s_1^{\text{task}} \times \mathbf{h}_{k1}^{\text{LM}} \\ \\ s_0^{\text{task}} \times \mathbf{h}_{k0}^{\text{LM}} \\ \quad\quad ([\mathbf{x}_k ; \mathbf{x}_k]) \end{cases}$$

Concatenate hidden layers

$[\overrightarrow{\mathbf{h}}_{kj}^{\text{LM}} ; \overleftarrow{\mathbf{h}}_{kj}^{\text{LM}}]$

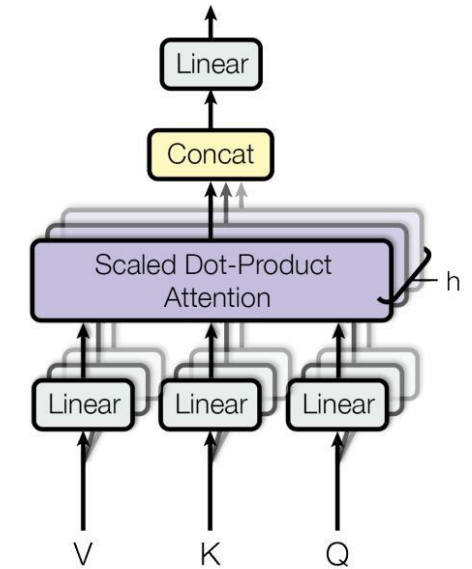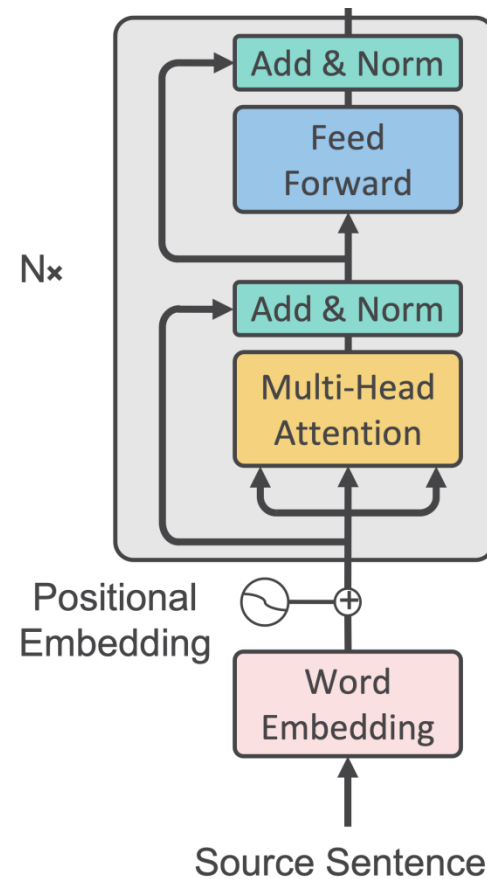Unlike usual word embeddings, ELMo is assigned to every *token* instead of a *type*

**biLMs**

Forward LM          Backward LM

$o_k$               $o_k$

$\overrightarrow{\mathbf{h}}_{k2}^{\text{LM}}$   $k-1$    $\overleftarrow{\mathbf{h}}_{k2}^{\text{LM}}$   $k+1$

$\overrightarrow{\mathbf{h}}_{k1}^{\text{LM}}$   $k-1$    $\overleftarrow{\mathbf{h}}_{k1}^{\text{LM}}$   $k+1$

$\mathbf{x}_k$

$t_k$               $t_k$

# Outline

- Introduction
- ELMO
- **BERT**

# Model Architecture

- **Multi-headed self attention**
  - Models context

- **Feed-forward layers**
  - Computes non-linear hierarchical features

- **Layer norm and residuals**
  - Makes training deep networks healthy

- **Positional embeddings**
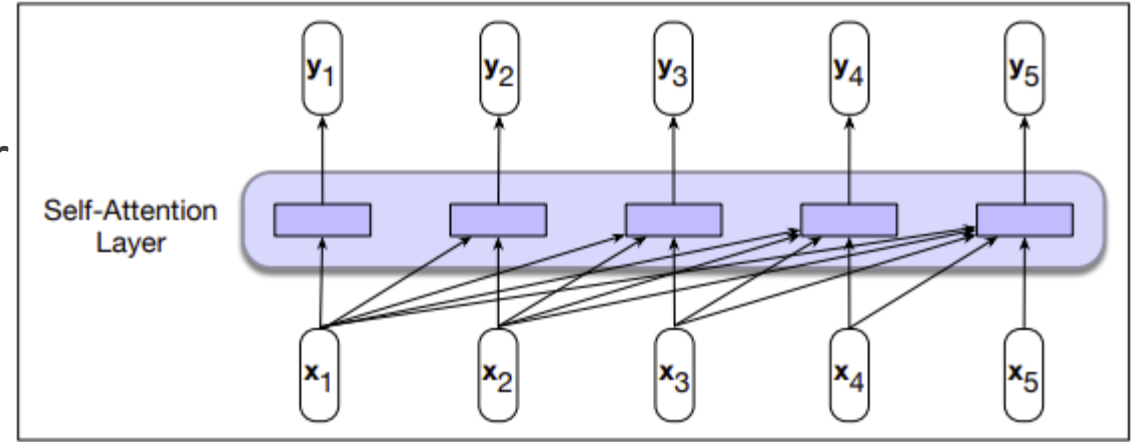  - Allows model to learn relative positioning

# Model Architecture

- Empirical advantages of Transformer vs. LSTM:

  ◦ Self-attention == no locality bias
    ◦ Long-distance context has "equal opportunity"

  ◦ Single multiplication per layer == efficiency on TPU
    ◦ Effective batch size is number of *words*, not *sequences*
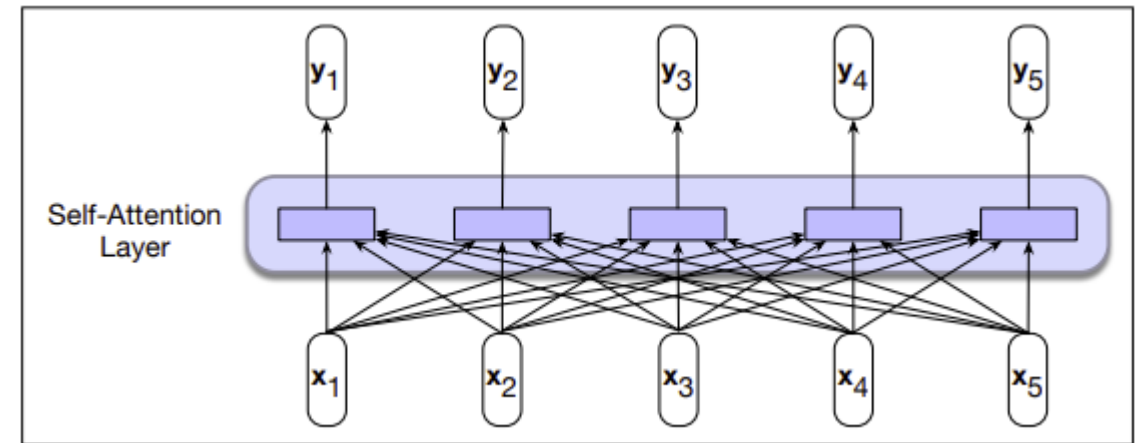
# Problem with Previous Methods

- **Problem**: Language models only use left context *or* right context, but language understanding is bidirectional.

- Why are LMs unidirectional?

- Reason 1: Directionality is needed to generate a well-formed probability distribution.
  - We don't care about this.

- Reason 2: Words can "see themselves" in a bidirectional encoder.
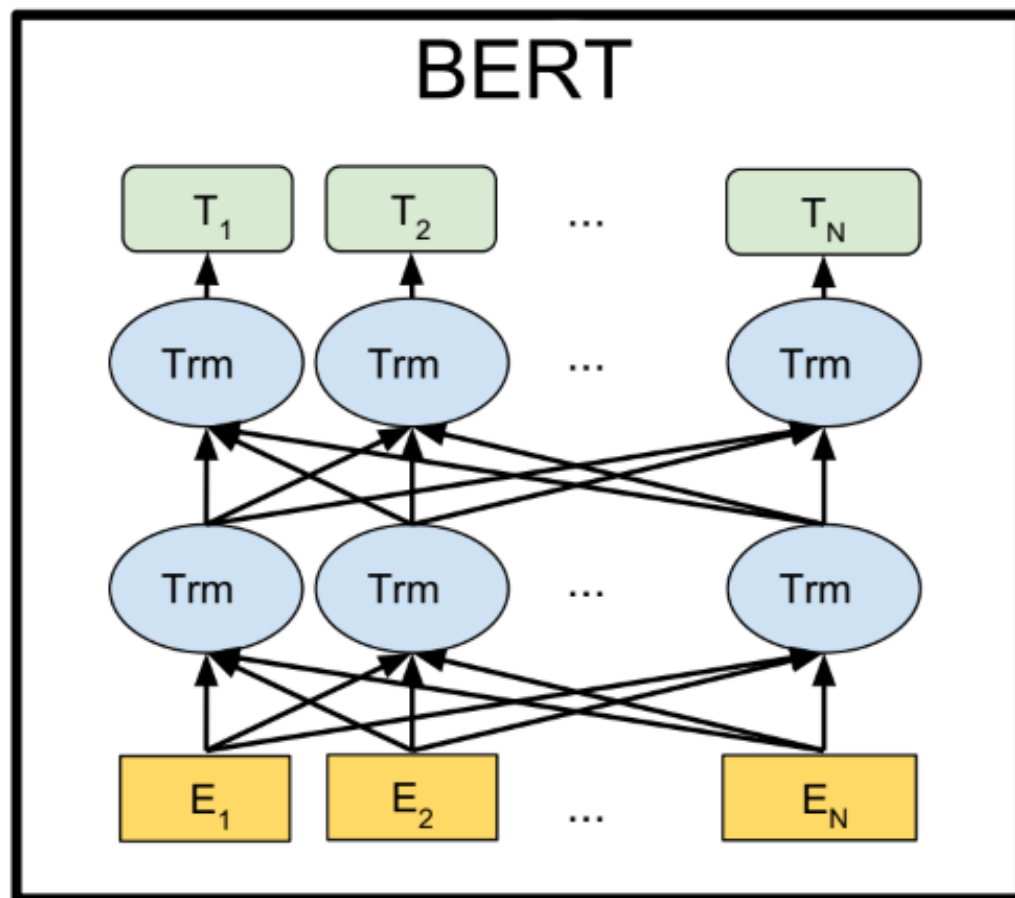
# Bidirectional Transformer

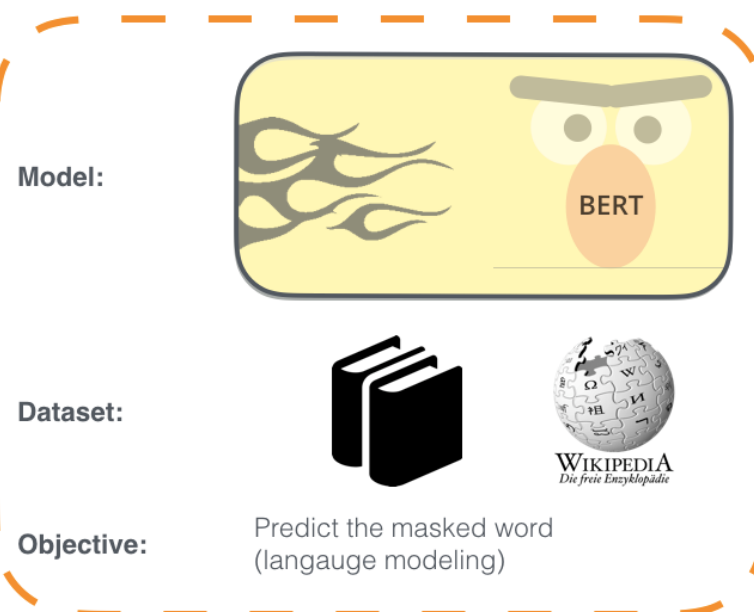- A causal backward looking transformer



- A bidirectional transformer



16

# BERT

# BERT

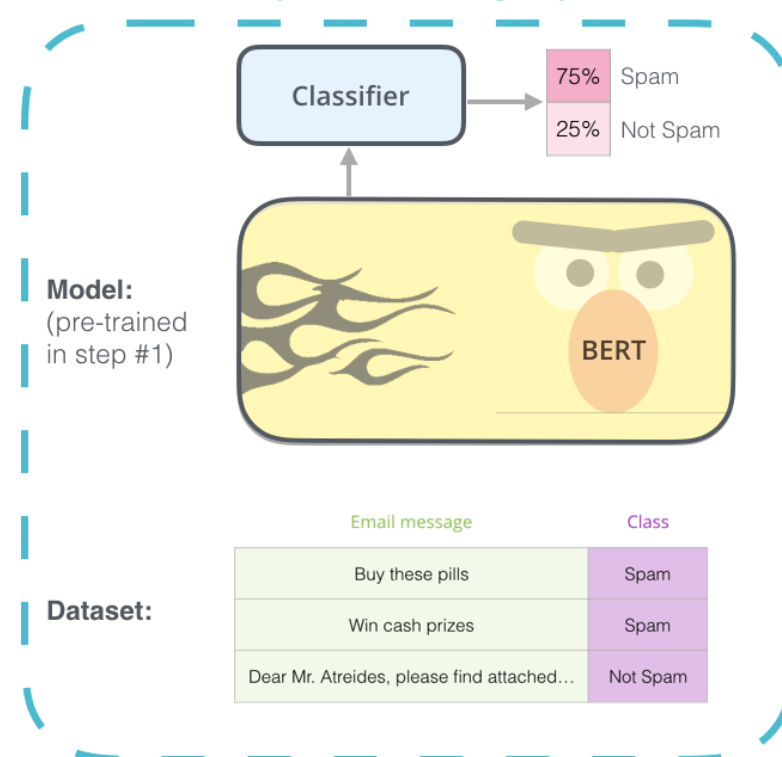1 - Semi-supervised training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.
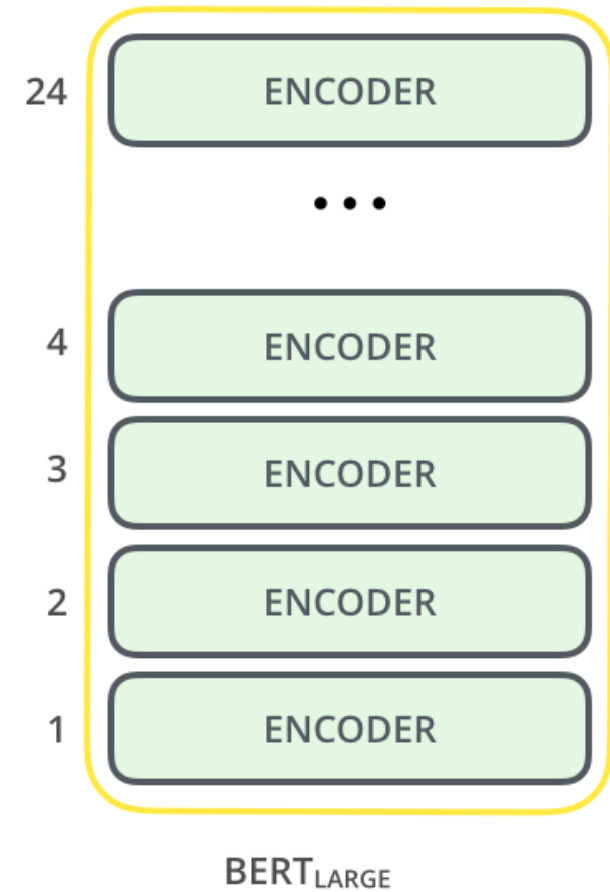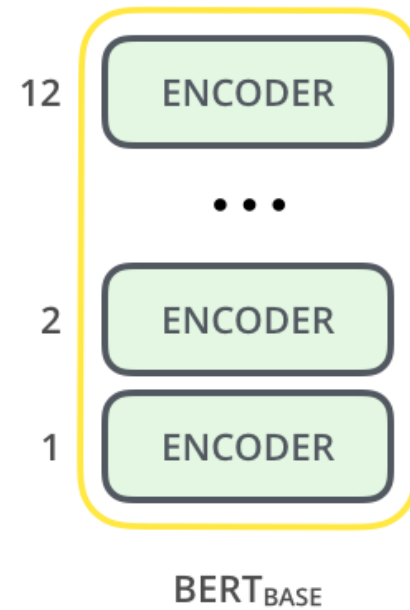
**Semi-supervised Learning Step**



**Model:**

**Dataset:**

**Objective:** Predict the masked word (langauge modeling)

2 - Supervised training on a specific task with a labeled dataset.

**Supervised Learning Step**



Classifier → 75% Spam / 25% Not Spam

**Model:** (pre-trained in step #1)

**Dataset:**

| Email message | Class |
|---|---|
| Buy these pills | Spam |
| Win cash prizes | Spam |
| Dear Mr. Atreides, please find attached… | Not Spam |

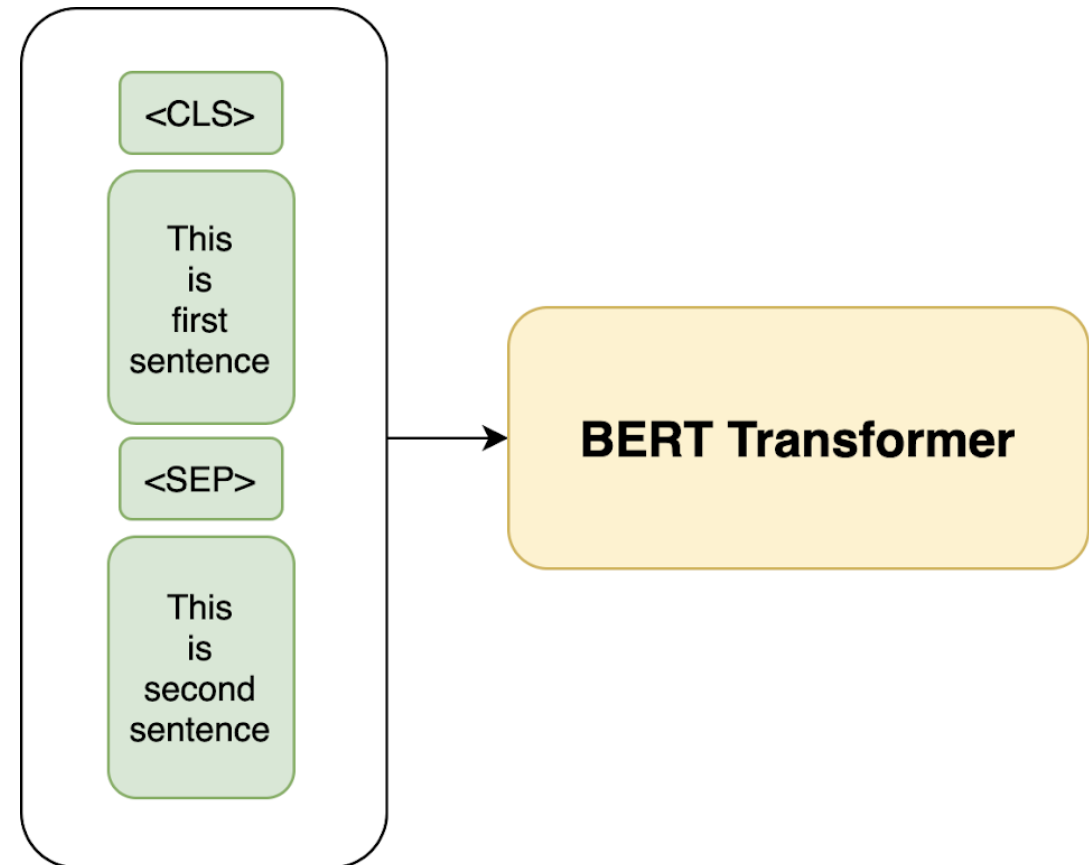# BERT

- BERT Base

- BERT Large



BERT_BASE                    BERT_LARGE
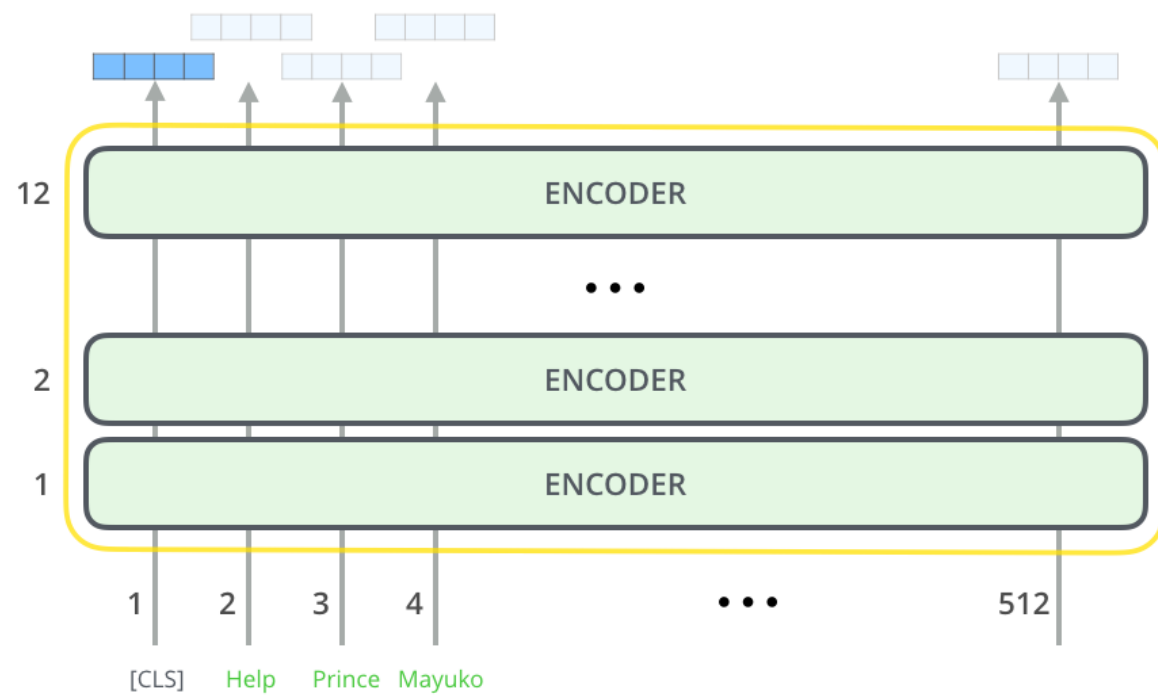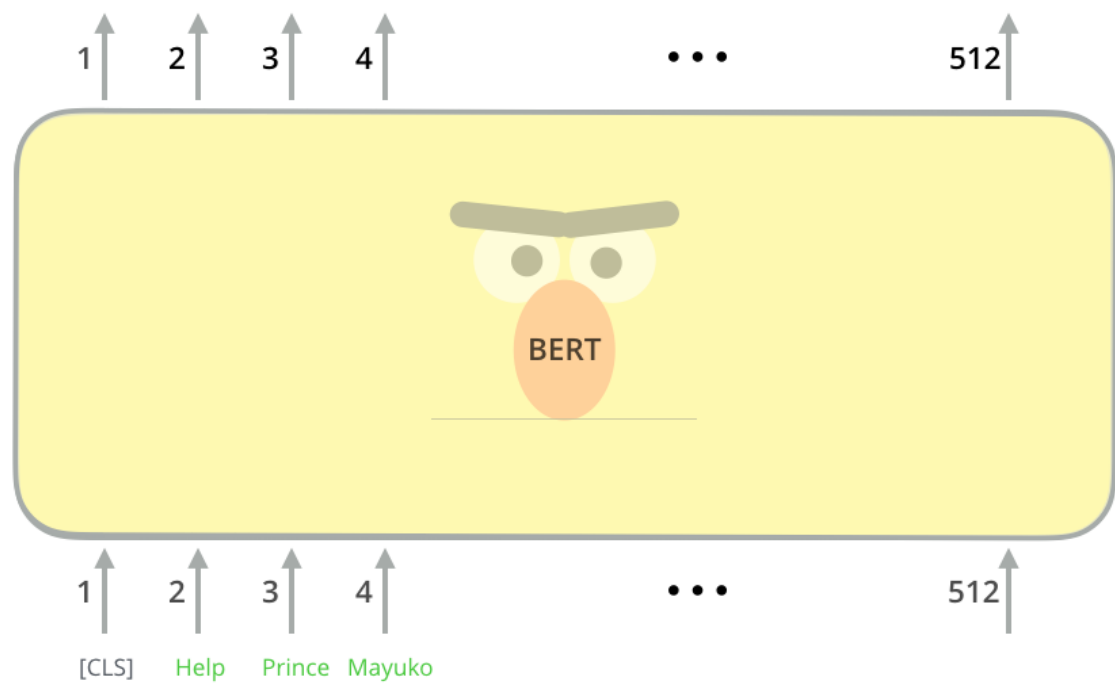
# BERT

- Training type
  - Masked language model
  - Next sentence prediction

- Input
  - First token  [CLS]
  - Delimiter token [SEP]
  - Masked token [MASK]

# BERT

# BERT

- Masked language model

- **Solution**: Mask out *k*% of the input words, and then predict the masked words
  - We always use *k* = 15%



- Too little masking: Too expensive to train

- Too much masking: Not enough context

# BERT

- Masked language model



Use the output of the masked word's position to predict the masked word

Possible classes: All English words

| 0.1% | Aardvark |
| ... | ... |
| 10% | Improvisation |
| ... | ... |
| 0% | Zyzzyva |

FFNN + Softmax

1  2  3  4  5  6  7  8  •••  512

BERT

Randomly mask 15% of tokens

1  2  3  4  5  6  7  8  •••  512

[CLS]  Let's  stick  to  [MASK]  in  this  skit

Input

[CLS]  Let's  stick  to improvisation in  this  skit

# BERT

- Next sentence prediction

- To learn *relationships* between sentences, predict whether Sentence B is actual sentence that proceeds Sentence A, or a random sentence

**Sentence A** = The man went to the store.
**Sentence B** = He bought a gallon of milk.
**Label** = IsNextSentence

**Sentence A** = The man went to the store.
**Sentence B** = Penguins are flightless.
**Label** = NotNextSentence
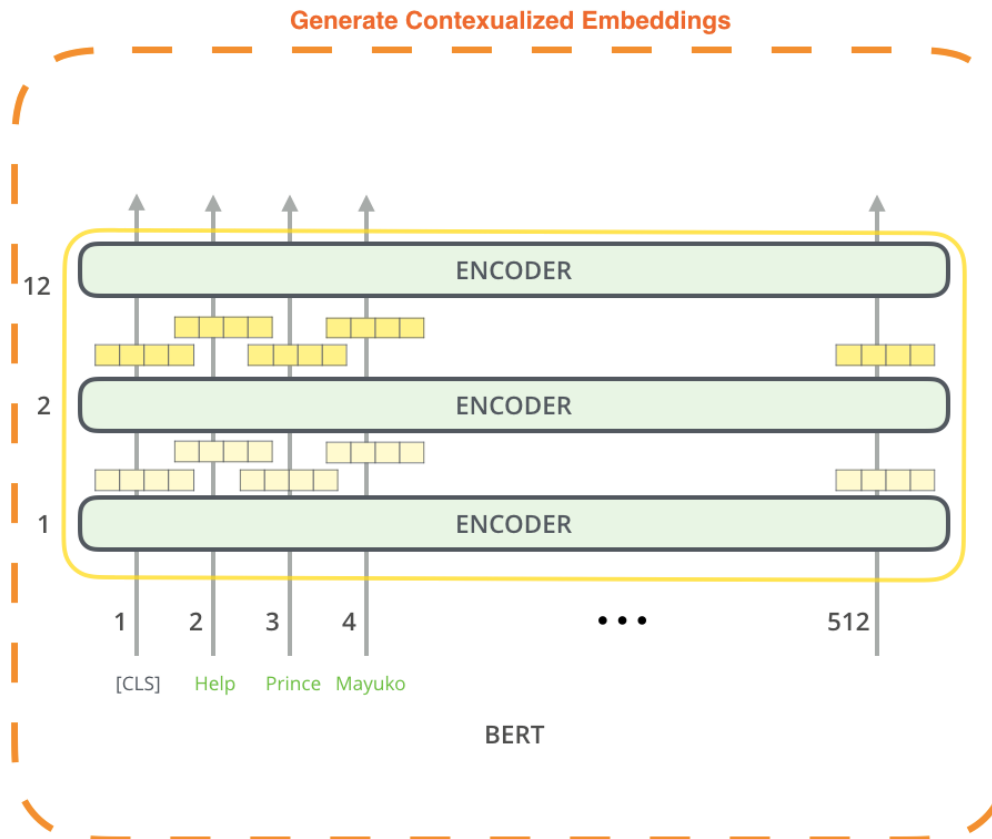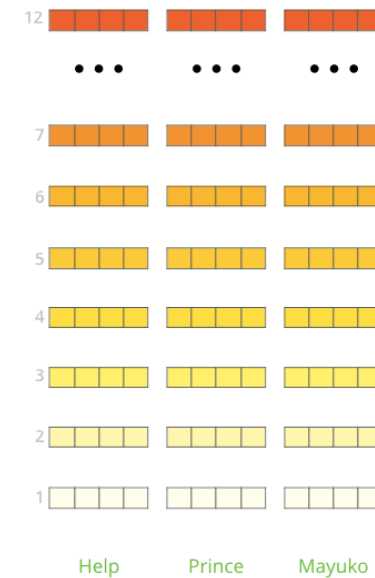
# BERT

- Next sentence prediction

Predict likelihood that sentence B belongs after sentence A

| 1% | IsNext |
| 99% | NotNext |

FFNN + Softmax

1  2  3  4  5  6  7  8  •••  512

BERT

Tokenized Input

1  2  •••  512

[CLS]  the  man  [MASK]  to  the  store  [SEP]

Input

[CLS] the man [MASK] to the store [SEP] penguin [MASK] are flightless birds [SEP]

Sentence A        Sentence B

# BERT



| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

# BERT

- Extracting embedding



**Generate Contexualized Embeddings**

ENCODER 12

ENCODER 2

ENCODER 1

1  2  3  4  · · ·  512

[CLS]  Help  Prince  Mayuko

BERT

The output of each encoder layer along each token's path can be used as a feature representing that token.

12

7

6

5

4

3

2

1

Help  Prince  Mayuko

**But which one should we use?**

# BERT

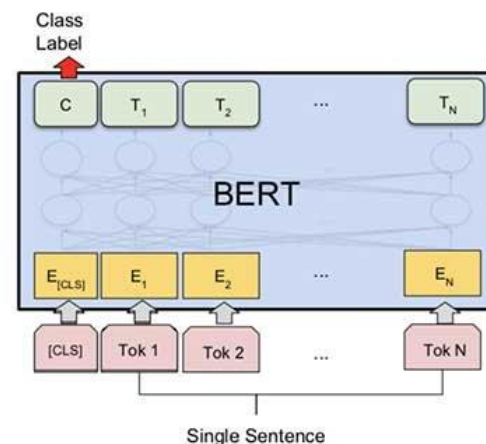- Extracting embedding

# Fine-Tuning Procedure
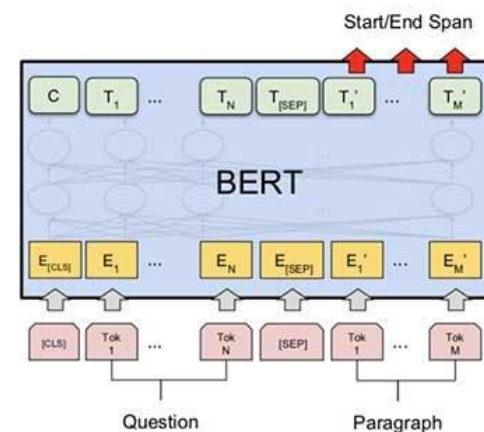
# Fine-Tuning Procedure

- Sentence pair classification

- Single sentence classification

- Question answering
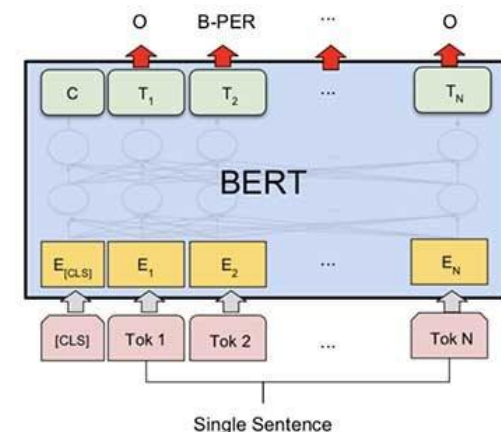
- Single sentence sequence labeling



(a) Sentence Pair Classification Tasks: MNLI, QQP, QNLI, STS-B, MRPC, RTE, SWAG

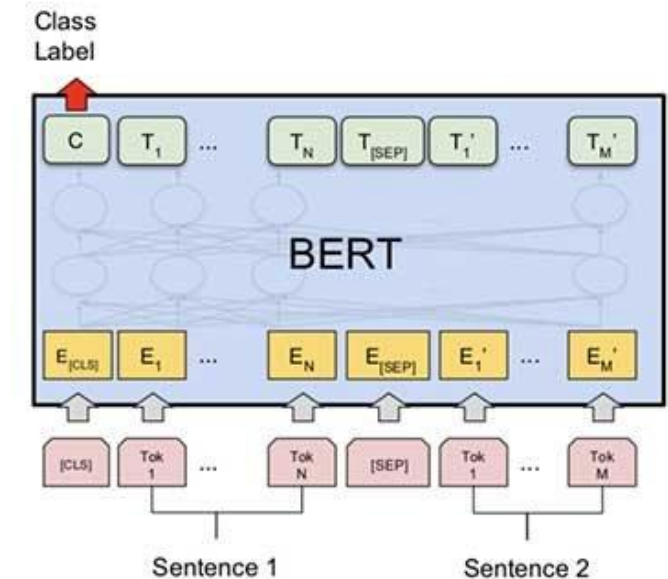(b) Single Sentence Classification Tasks: SST-2, CoLA

(c) Question Answering Tasks: SQuAD v1.1

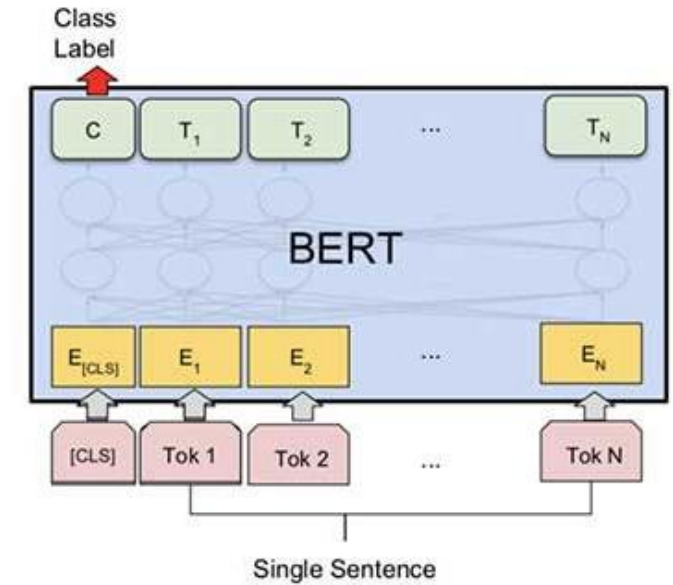(d) Single Sentence Tagging Tasks: CoNLL-2003 NER

# Fine-Tuning Procedure

- Sentence pair classification tasks
  - Paraphrase identification
  - Answer retrieval
  - Textual entailment



- Datasets:
  - MNLI
  - QQP
  - QNLI
  - STS-B
  - MRPC
  - RTE
  - SWAG

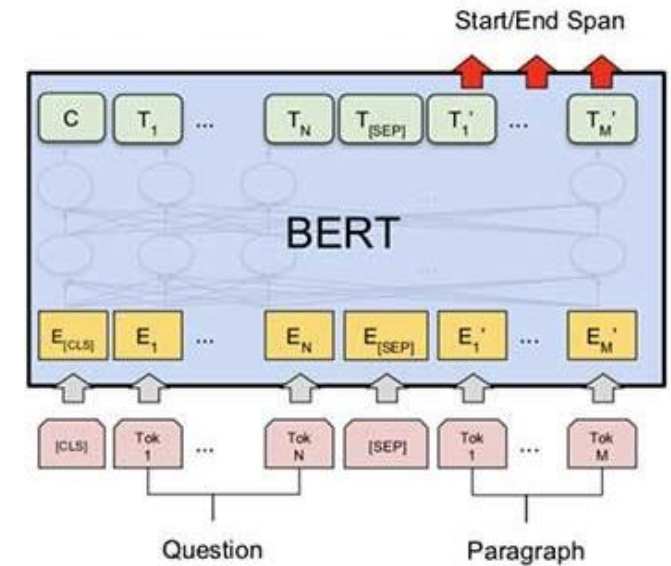# Fine-Tuning Procedure

- Single sentence classification tasks:
  - Spam detection
  - Sentiment analysis
  - News categorization


- Datasets:
  - SST-2
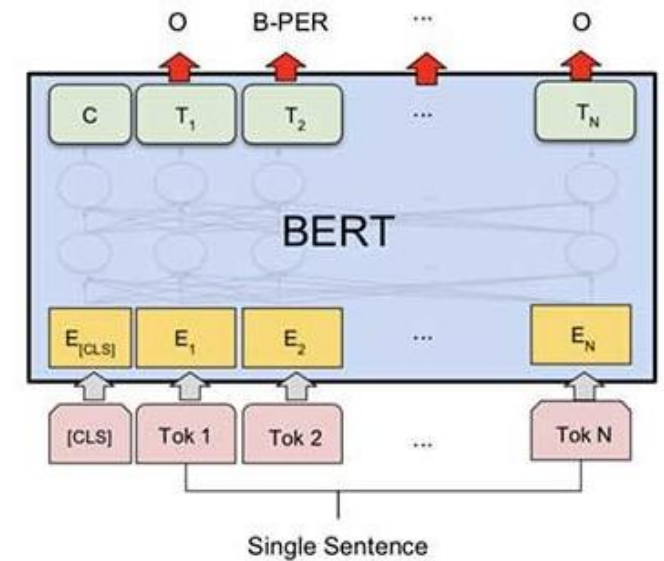  - CoLA

# Fine-Tuning Procedure

- Question answering tasks
  (finding answer span)

- Dataset:
  ◦ SQuAD

# Fine-Tuning Procedure

- Single sentence sequence labeling tasks:
  - NER
  - POS tagging
  - Slot filing

- Datasets:
  - CoNLL-2003 NER

# Motivations for Improving BERT

- Accuracy

- Large

- Slow

- Hard to train

| | **BERT** |
|---|---|
| **Size (millions)** | **Base**: 110 <br> **Large**: 340 |
| **Training Time** | **Base**: 8 x V100 x 12 days* <br> **Large:** 64 TPU Chips x 4 days (or 280 x V100 x 1 days*) |
| **Performance** | Outperforms state-of-the-art in Oct 2018 |
| **Data** | 16 GB BERT data (Books Corpus + Wikipedia). <br> 3.3 Billion words. |
| **Method** | BERT (Bidirectional Transformer with MLM and NSP) |

# Model Improvement Methods

- Larger with more data

- Quantization
  ◦ Quantization-aware training

- Distillation

- Pruning

- Specialization

# Further Reading

- Speech and Language Processing (3$^{rd}$ ed. draft)
  - Chapter 11