# Tokenisation for Transformers & Contextual Embeddings

## Introduction to NLP

**Rahul Mishra**

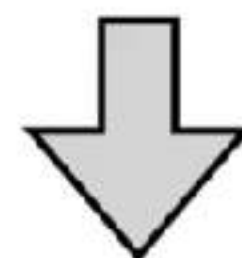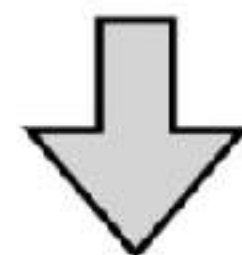IIIT-Hyderabad

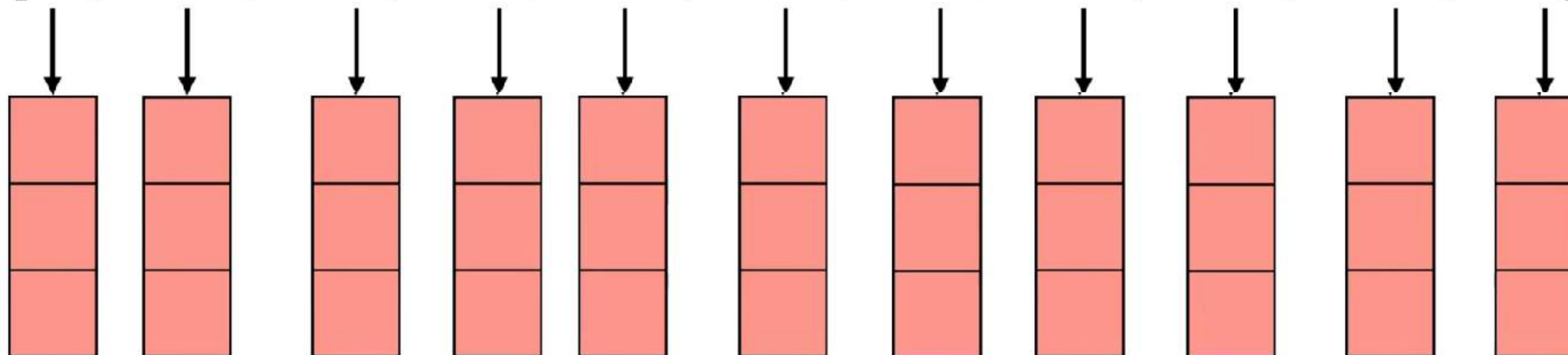March 22, 2024

# Byte Pair Encoding

"Today is a beautiful day outside."

⬇

["To", "day", "is", "a", "beaut", "iful", "day", "out", "side", "."]

⬇

[98, 1452,    43,   15, 2932, 1709,  740, 1452, 3112, 3823,  74]
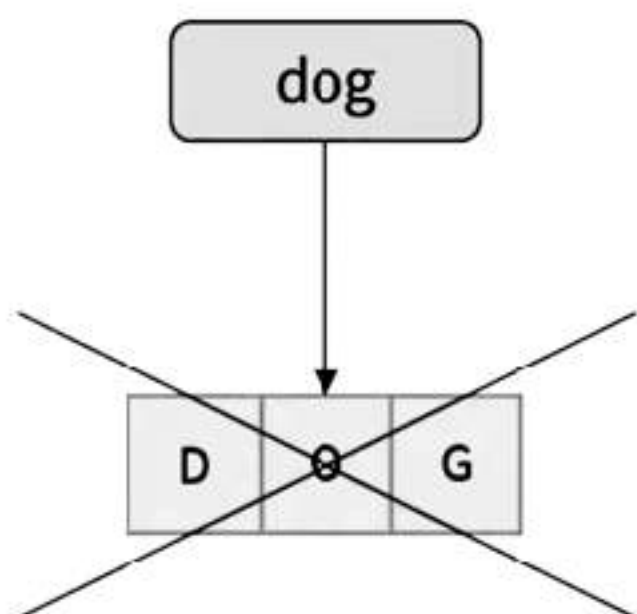
# Byte Pair Encoding

Instead of

- white-space segmentation
- single-character segmentation

**Use the data** to tell us how to tokenize.

**Subword tokenization** (because tokens can be parts of words as well as whole words)

# Byte Pair Encoding

Frequently used words should not be split into smaller subwords

Rare words should be decomposed into meaningful subwords.

Word-based tokenization

Character-based tokenization

dog

D O G

dogs

dog | s

Very large vocabularies

Large quantity of out-of-vocabulary tokens

Loss of meaning across very similar words

Very long sequences

Less meaningful individual tokens

# Byte Pair Encoding

Three common algorithms:
- **Byte-Pair Encoding (BPE)** (Sennrich et al., 2016)
- **Unigram language modeling tokenization** (Kudo, 2018)
- **WordPiece** (Schuster and Nakajima, 2012)

All have 2 parts:
- A token **learner** that takes a raw training corpus and induces a vocabulary (a set of tokens).
- A token **segmenter** that takes a raw test sentence and tokenizes it according to that vocabulary

# Byte Pair Encoding

Let vocabulary be the set of all individual characters
= {A, B, C, D,…, a, b, c, d….}

Repeat:

◦ Choose the two symbols that are most frequently adjacent in the training corpus (say 'A', 'B')

◦ Add a new merged symbol 'AB' to the vocabulary

◦ Replace every adjacent 'A' 'B' in the corpus with 'AB'.

Until $k$ merges have been done.

# BPE Token Learning Algo

**function** BYTE-PAIR ENCODING(strings $C$, number of merges $k$) **returns** vocab $V$

$V \leftarrow$ all unique characters in $C$        # initial set of tokens is characters
**for** $i = 1$ **to** $k$ **do**        # merge tokens til $k$ times
    $t_L, t_R \leftarrow$ Most frequent pair of adjacent tokens in $C$
    $t_{NEW} \leftarrow t_L + t_R$        # make new token by concatenating
    $V \leftarrow V + t_{NEW}$        # update the vocabulary
    Replace each occurrence of $t_L, t_R$ in $C$ with $t_{NEW}$        # and update the corpus
**return** $V$

# BPE Token Learner

Original (very fascinating🙄) corpus:

low low low low low lowest lowest newer newer newer
newer newer newer wider wider wider new new

Add end-of-word tokens, resulting in this vocabulary:

**vocabulary**
_, d, e, i, l, n, o, r, s, t, w

**corpus representation**

| | |
|---|---|
| 5 | l o w _ |
| 2 | l o w e s t _ |
| 6 | n e w e r _ |
| 3 | w i d e r _ |
| 2 | n e w _ |

# BPE Token Learner

| corpus | | vocabulary |
|---|---|---|
| 5 | l o w _ | _, d, e, i, l, n, o, r, s, t, w |
| 2 | l o w e s t _ | |
| 6 | n e w e r _ | |
| 3 | w i d e r _ | |
| 2 | n e w _ | |

## Merge e r to er

| corpus | | vocabulary |
|---|---|---|
| 5 | l o w _ | _, d, e, i, l, n, o, r, s, t, w, er |
| 2 | l o w e s t _ | |
| 6 | n e w er _ | |
| 3 | w i d er _ | |
| 2 | n e w _ | |

# BPE Token Learner

**corpus**

| | |
|---|---|
| 5 | l o w _ |
| 2 | l o w e s t _ |
| 6 | n e w er_ |
| 3 | w i d er_ |
| 2 | n e w _ |

**vocabulary**

_, d, e, i, l, n, o, r, s, t, w, er, er_

## Merge n e to ne

**corpus**

| | |
|---|---|
| 5 | l o w _ |
| 2 | l o w e s t _ |
| 6 | ne w er_ |
| 3 | w i d er_ |
| 2 | ne w _ |

**vocabulary**

_, d, e, i, l, n, o, r, s, t, w, er, er_, ne

# BPE Token Learner

The next merges are:

| Merge | Current Vocabulary |
|---|---|
| (ne, w) | _, d, e, i, l, n, o, r, s, t, w, er, er_, ne, new |
| (l, o) | _, d, e, i, l, n, o, r, s, t, w, er, er_, ne, new, lo |
| (lo, w) | _, d, e, i, l, n, o, r, s, t, w, er, er_, ne, new, lo, low |
| (new, er_) | _, d, e, i, l, n, o, r, s, t, w, er, er_, ne, new, lo, low, newer_ |
| (low, _) | _, d, e, i, l, n, o, r, s, t, w, er, er_, ne, new, lo, low, newer_, low_ |

# BPE Token Segmenter

On the test data, run each merge learned from the training data:

- Greedily
- In the order we learned them
- (test frequencies don't play a role)

So: merge every e r to er, then merge er _ to er_, etc.

Result:

- Test set "n e w e r _" would be tokenized as a full word
- Test set "l o w e r _" would be two tokens: "low er_"

# BPE Properties

Usually include frequent words

And frequent subwords

- Which are often morphemes like *-est* or *–er*

A **morpheme** is the smallest meaning-bearing unit of a language

- *unlikeliest* has 3 morphemes *un-*, *likely*, and *-est*

# WordPiece and SentencePiece

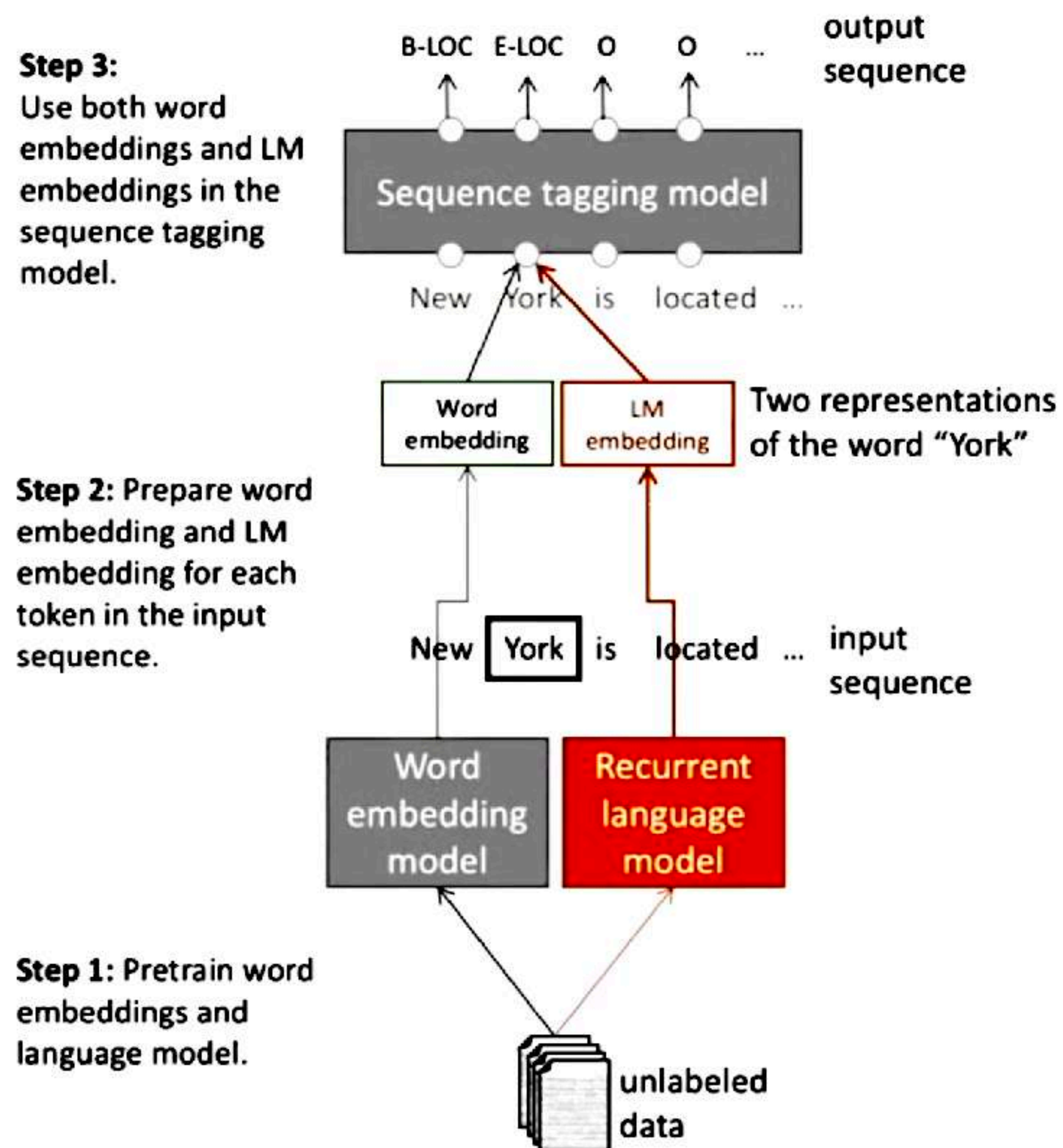$$\frac{P(tok_1, tok_2)}{P(tok_1)P(tok_2)} \longrightarrow \text{Pair probability}$$
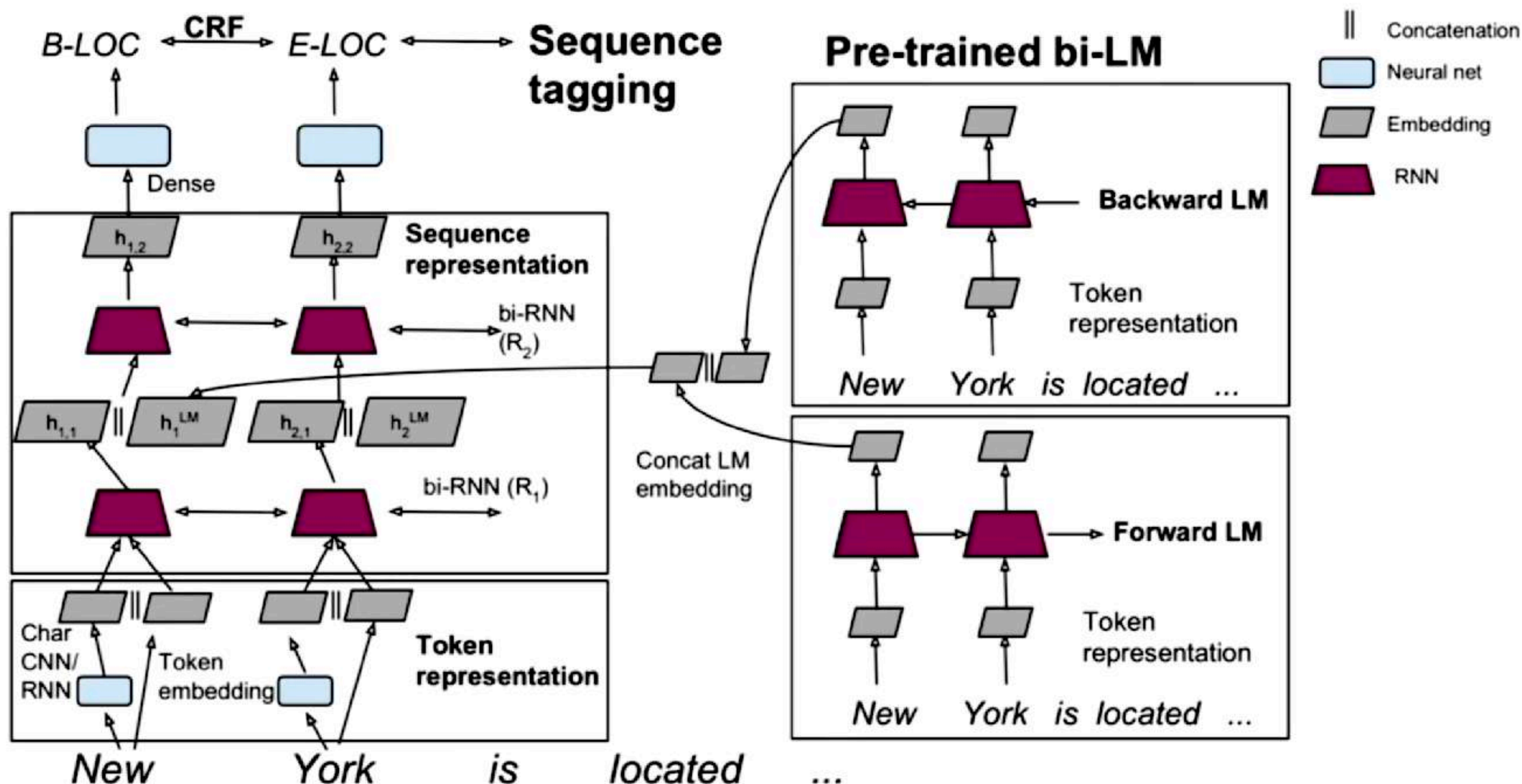
**Japanese**

トークナイザーに
ついて楽しく学ん
でいただければ幸
いです

**Chinese**

我希望您在学习
分词器过程中度
过愉快的时光

How do you pretokenize such languages?

# TagLM (Pre-ELMo) Peters et al. 2017



**Step 3:** Use both word embeddings and LM embeddings in the sequence tagging model.

**Step 2:** Prepare word embedding and LM embedding for each token in the input sequence.

**Step 1:** Pretrain word embeddings and language model.

B-LOC  E-LOC  O  O  ...  output sequence

Sequence tagging model

New  York  is  located  ...

Word embedding | LM embedding — Two representations of the word "York"

New  York  is  located  ...  input sequence

Word embedding model | Recurrent language model

unlabeled data

# TagLM (Pre-ELMo) Peters et al. 2017



$$\mathbf{h}_{k,1} = [\overrightarrow{\mathbf{h}}_{k,1}; \overleftarrow{\mathbf{h}}_{k,1}; \mathbf{h}_k^{LM}].$$

# ELMo: Embeddings from Language Models

Train a bidirectional LM

Aim at performant but not overly large LM:

- Use 2 biLSTM layers
- Use character CNN to build initial word representation (only)
  - 2048 char n-gram filters and 2 highway layers, 512 dim projection
- User 4096 dim hidden/cell LSTM states with 512 dim projections to next input
- Use a residual connection
- Tie parameters of token input and output (softmax) and tie these between forward and backward LMs

Best Paper award at NAACL 2018

# ELMo: Embeddings from Language Models

- ELMo learns task-specific combination of biLM representations
- This is an innovation that improves on just using top layer of LSTM stack

$$R_k = \{\mathbf{x}_k^{LM}, \overrightarrow{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \ldots, L\}$$

$$= \{\mathbf{h}_{k,j}^{LM} \mid j = 0, \ldots, L\},$$

$$\mathbf{ELMo}_k^{task} = E(R_k; \Theta^{task}) = \gamma^{task} \sum_{j=0}^{L} s_j^{task} \mathbf{h}_{k,j}^{LM}$$

- $\gamma^{task}$ scales overall usefulness of ELMo to task;
- $\mathbf{s}^{task}$ are softmax-normalized mixture model weights

# ELMo

| ELMo | ELMo in BiLSTM | 2018 | 92.22 |
|------|----------------|------|-------|
| TagLM Peters | LSTM BiLM in BiLSTM tagger | 2017 | 91.93 |
| Ma + Hovy | BiLSTM + char CNN + CRF layer | 2016 | 91.21 |
| Tagger Peters | BiLSTM + char CNN + CRF layer | 2017 | 90.87 |
| Ratinov + Roth | Categorical CRF+Wikipeda+word cls | 2009 | 90.80 |
| Finkel et al. | Categorical feature CRF | 2005 | 86.86 |
| IBM Florian | Linear/softmax/TBL/HMM ensemble, gazettes++ | 2003 | 88.76 |
| Stanford | MEMM softmax markov model | 2003 | 86.07 |

| TASK | PREVIOUS SOTA | | OUR BASELINE | ELMo + BASELINE | INCREASE (ABSOLUTE/ RELATIVE) |
|------|---------------|---|--------------|-----------------|-------------------------------|
| SQuAD | Liu et al. (2017) | 84.4 | 81.1 | 85.8 | 4.7 / 24.9% |
| SNLI | Chen et al. (2017) | 88.6 | 88.0 | $88.7 \pm 0.17$ | 0.7 / 5.8% |
| SRL | He et al. (2017) | 81.7 | 81.4 | 84.6 | 3.2 / 17.2% |
| Coref | Lee et al. (2017) | 67.2 | 67.2 | 70.4 | 3.2 / 9.8% |
| NER | Peters et al. (2017) | $91.93 \pm 0.19$ | 90.15 | $92.22 \pm 0.10$ | 2.06 / 21% |
| SST-5 | McCann et al. (2017) | 53.7 | 51.4 | $54.7 \pm 0.5$ | 3.3 / 6.8% |

# ELMo layers

The two biLSTM NLM layers have differentiated uses/meanings

- Lower layer is better for lower-level syntax, etc.
  - Part-of-speech tagging, syntactic dependencies, NER
- Higher layer is better for higher-level semantics
  - Sentiment, Semantic role labeling, question answering, SNLI

# ULMfit



ULMfit
Jan 2018
Training:
1 GPU day

GPT
June 2018
Training
240 GPU days

BERT
Oct 2018
Training
256 TPU days
~320–560
GPU days

GPT-2
Feb 2019
Training
~2048 TPU v3
days according to
a reddit thread

# BERT

**Problem**: Language models only use left context *or* right context, but language understanding is bidirectional.
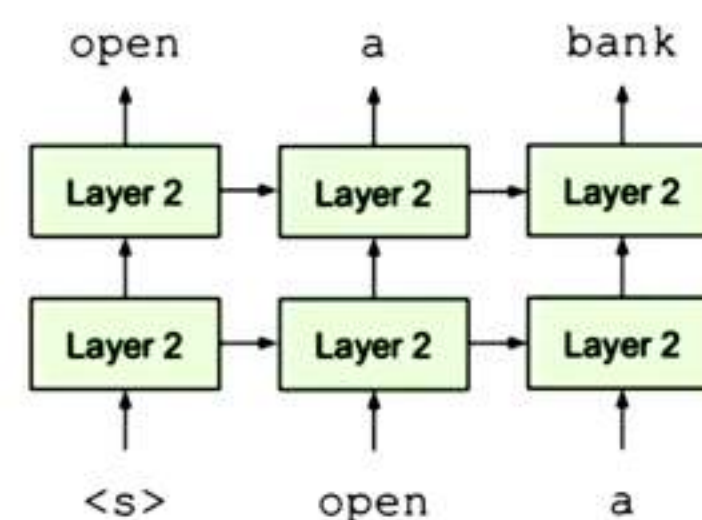
Why are LMs unidirectional?

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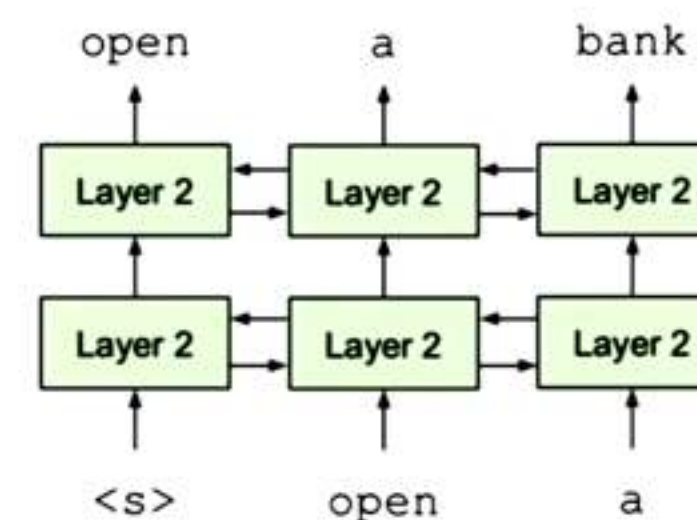
- We don't care about this.

Reason 2: Words can "see themselves" in a bidirectional encoder.



**Unidirectional context**
Build representation incrementally

**Bidirectional context**
Words can "see themselves"

# BERT

**Solution**: Mask out $k$% of the input words, and then predict the masked words

- They always use $k$ = 15%
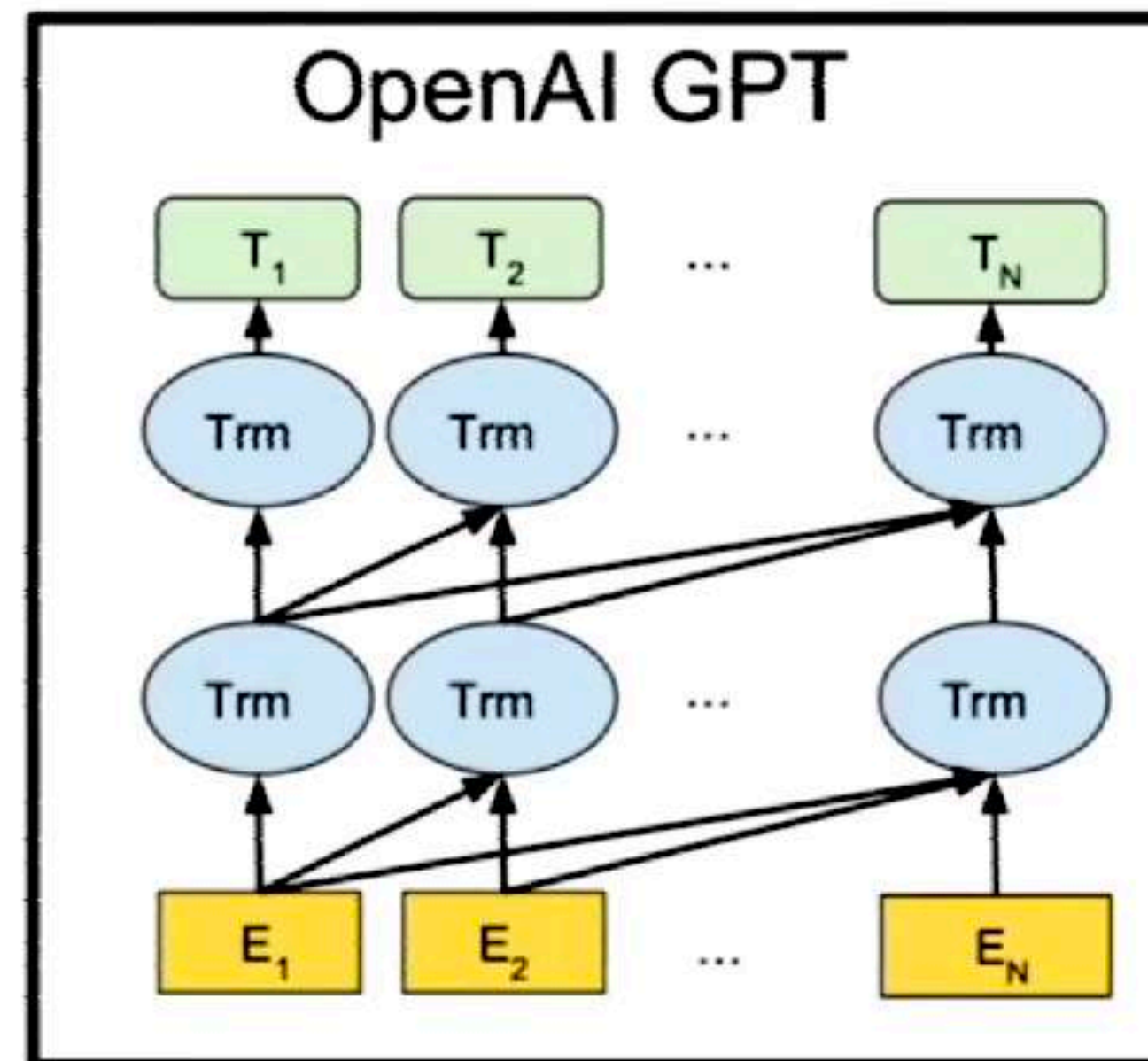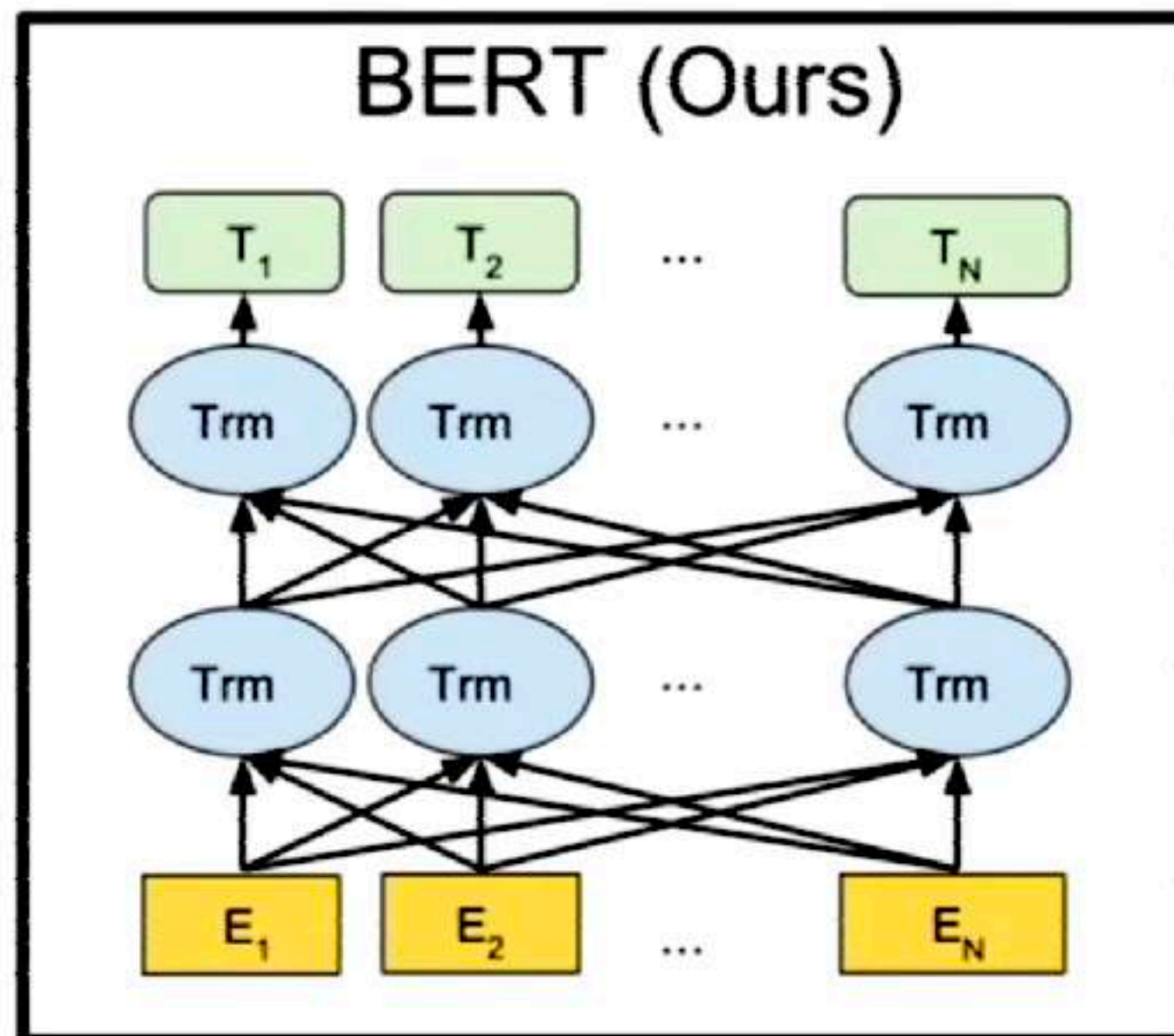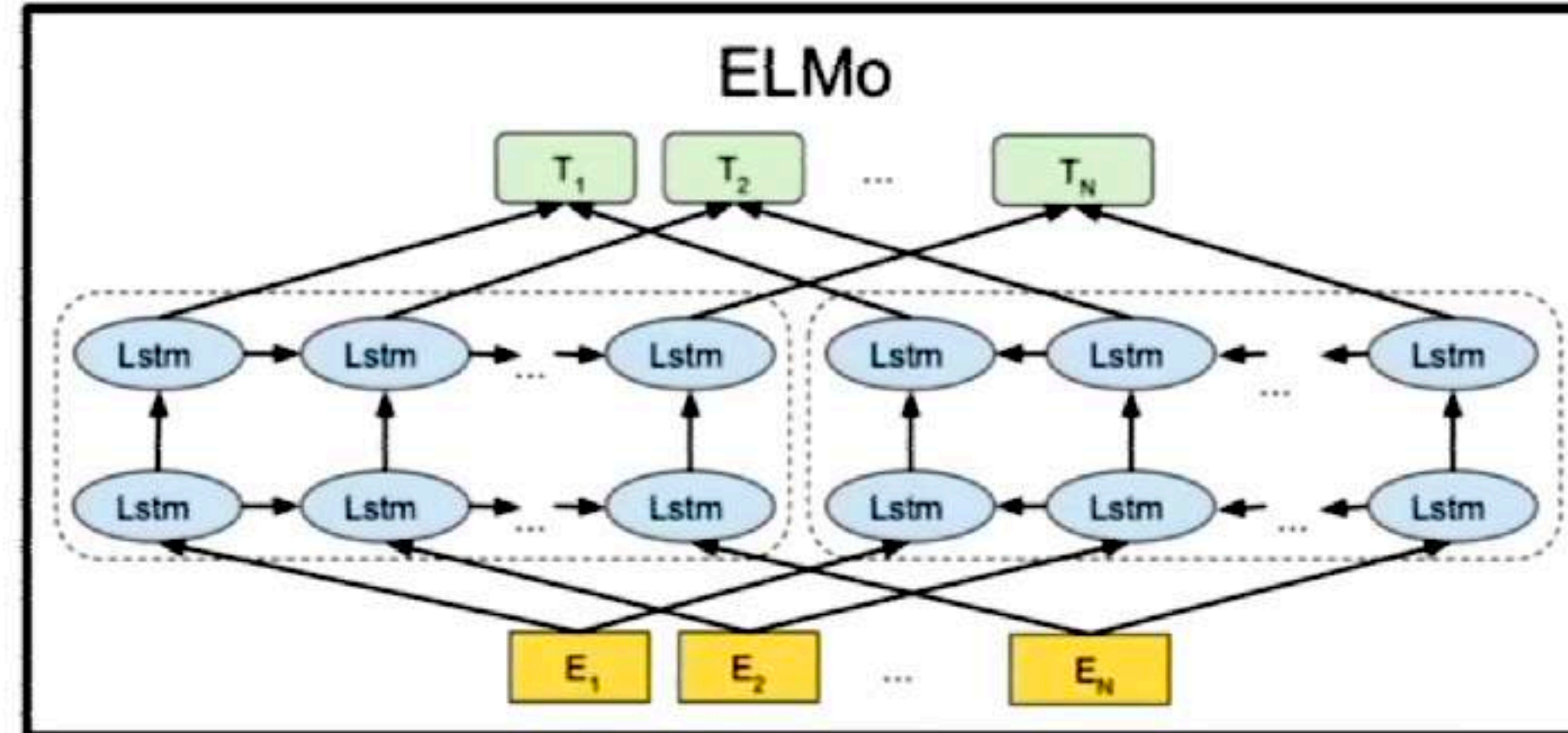
<pre>
                    store           gallon
                     ↑                ↑
   the man went to the [MASK] to buy a [MASK] of milk
</pre>

Too little masking: Too expensive to train

Too much masking: Not enough context

# BERT

# BERT Fine Tuning

- Simply learn a classifier built on the top layer for each task that you fine tune for

# 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 Sent Pair Encoding



Token embeddings are word pieces
Learned segmented embedding represents each sentence
Positional embedding is as for other Transformer architectures

A timeline of large language models (LLMs) from 2019 to 2023.

**2019**
- G T5
- GPT-3

**2020**
- G GShard

**2021**
  - **1-4**
    - G mT5
    - HUAWEI PanGu-α
    - PLUG
    - Open-Source
    - AI21 labs Jurassic-1
    - NAVER LABS HyperCLOVA
    - G LaMDA
    - BAAI CPM-2
  - **5-8**
    - Ernie 3.0
    - BLOOM
    - Codex
  - **9-10**
    - G FLAN
    - inspur Yuan 1.0
    - MT-NLG
    - T0
  - **11-12**
    - Gopher
    - G GLaM

**2022**
- WebGPT
- Ernie 3.0 Titan
  - **1-3**
    - InstructGPT
    - GPT-NeoX-20B
    - CodeGen
    - G UL2
    - G PaLM
  - **4-6**
    - Tk-Instruct AI2
    - OPT
    - mT0
    - BLOOMZ
    - Galatica
    - OPT-IML
    - GLM
    - AlexaTM
  - **7-10**
    - Sparrow
    - G Flan-T5
    - G Flan-PaLM
    - AlphaCode
    - Chinchilla
    - HUAWEI PanGu-Σ
  - **11-12**
    - ChatGPT

**2023**
  - **1-3**
    - G Bard
    - ERNIE Bot
    - LLaMA
    - GPT-4