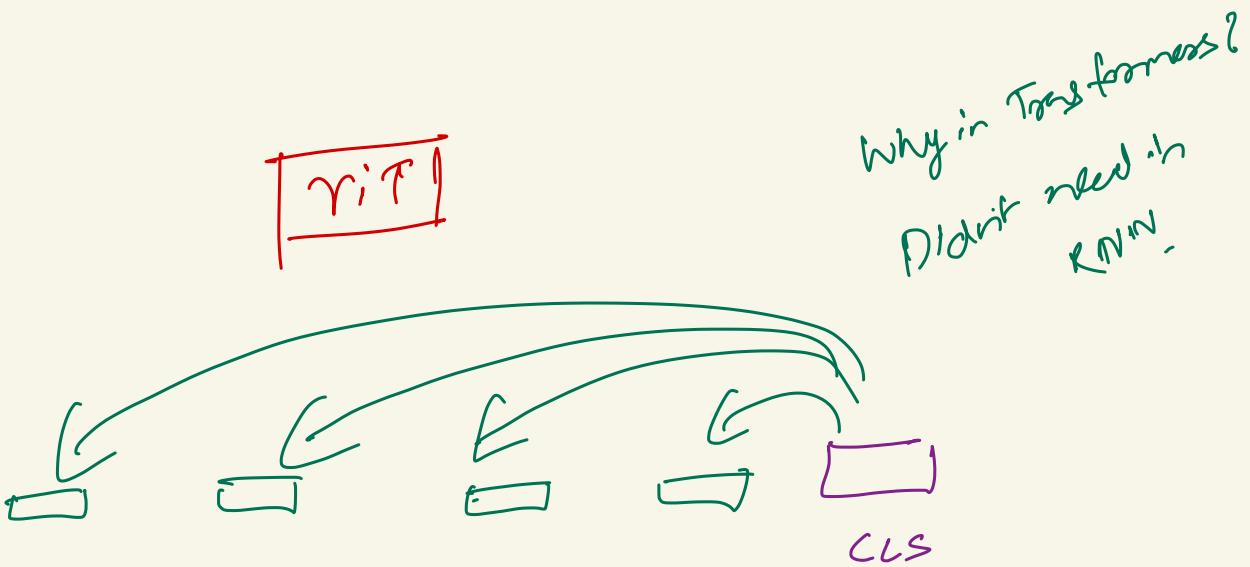
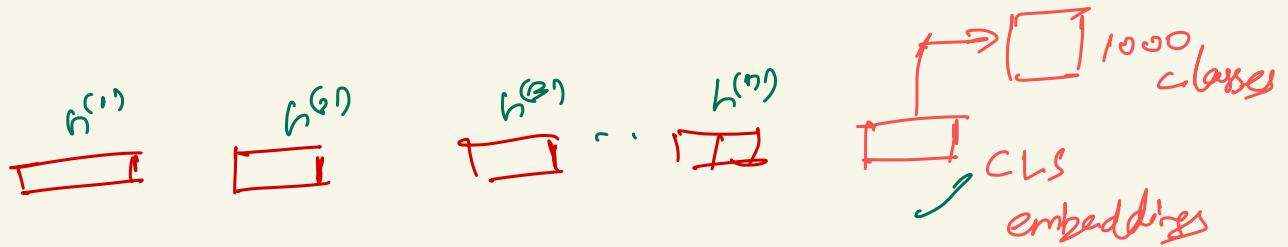


More Details about BERT

$(12 \times 12^f \times (768)^2) / 85m$

- Two models were released:
 - BERT-base: 12 layers, 768-dim hidden states, 12 attention heads, 110 million params.
 - BERT-large: 24 layers, 1024-dim hidden states, 16 attention heads, 340 million params.
- Trained on:
 - BooksCorpus (800 million words)
 - English Wikipedia (2,500 million words)
- Pretraining is expensive and impractical on a single GPU.
 - BERT was pretrained with 64 TPU chips for a total of 4 days.
 $30k \times 768 \approx 2.5m$
 - (TPUs are special tensor operation acceleration hardware)
- Finetuning is practical and common on a single GPU
 - “Pretrain once, finetune many times.”

BERT vocabulary: 30,522



Encoder vs Decoder

Same # of parameters

Same modules

Bi-directional attn

Causal attention

Classification

Generation

Sig. labeling

Sentence rep. (2 sentences are
similar?)

Essential for Pretraining /

1. Data (lots & lots of)

?, Compute

Model Architecture

self-supervision

BERT uses subwords

GPT-3.5 & GPT-4 GPT-3 (Legacy)

Have the bards who preceded me left any theme unsung?

Clear

Show example

Tokens

13

Characters

53

Have the bards who preceded me left any theme unsung?

Text

Token IDs

The original word 'embeddings' is denoted as:
['em', '#bed', '#ding', '#s']

Byte-Pair-Encoding (BPE)

[coined by [Gage et al., 1994](#); adapted to the task of word segmentation by [Sennrich et al., 2016](#); see [Gallé \(2019\)](#) for more]

Main idea: Use our data to automatically tell us what the tokens should be

Token learner:

Raw train corpus \Rightarrow Vocabulary (a set of tokens)

} **BPE**

Token segmenter:

Raw sentences \Rightarrow Tokens in the vocabulary

Types are distinct tokens in a corpus

Corpus vs. Dataset

Byte-Pair-Encoding (BPE) – Token learner

[coined by [Gage et al., 1994](#); adapted to the task of word segmentation by [Sennrich et al., 2016](#); see [Gallé \(2019\)](#) for more]

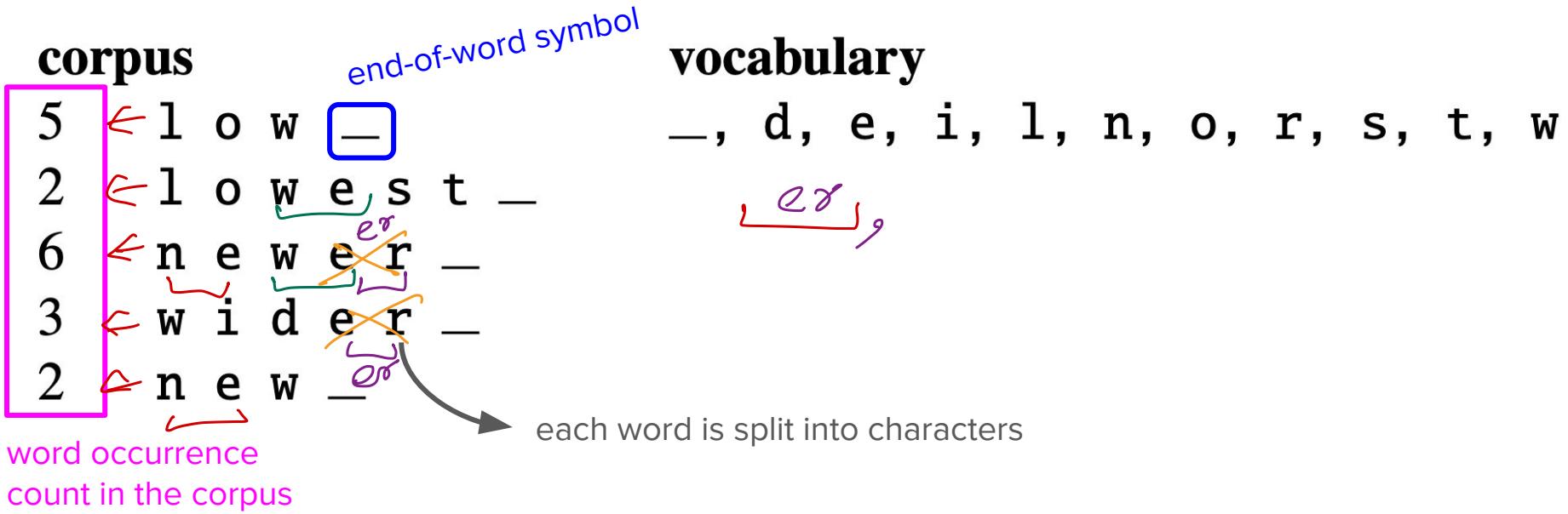
Raw train corpus \Rightarrow Vocabulary (a set of tokens)

- ⌚ Pre-tokenize the corpus in words & append a special end-of-word symbol E to each word
- ⌚ Initialize vocabulary with the set of all individual characters
- ⌚ Choose 2 tokens that are most frequently adjacent (“A”, “B”)
 - ⌚ Respect word boundaries: Run the algorithm inside words
- ⌚ Add a new merged symbol (“AB”) to the vocabulary
- ⌚ Change the occurrence of the 2 selected tokens with the new merged token in the corpus
- ⌚ Continues doing this until k merges are done

All k new symbols and initial characters are the final vocabulary

What's k ? Open research question, see [Mielke et al., 2021](#) Sec 6.6; 30K is seen frequently

Byte-Pair-Encoding (BPE) – Token learner *Example*



Byte-Pair-Encoding (BPE) – Token learner *Example*

- Counts all pairs of adjacent symbols
- The most frequent is the pair **e r** [a total of 9 occurrences]
- Merge these symbols, treating **er** as one symbol, & add the new symbol to the vocabulary

corpus

5	l	o	w	_			
2	l	o	w	e	s	t	_
6	n	e	w	er	L		
3	w	i	d	er	L		
2	n	e	w	_			

vocabulary

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

e, r → er
er, — → er-

Byte-Pair-Encoding (BPE) – Token learner *Example*

- Counts all pairs of adjacent symbols
- The most frequent is the pair **er_**
- Merge these symbols, treating **er_** as one symbol, & add the new symbol to the vocabulary

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**
me, new,

Byte-Pair-Encoding (BPE) – Token learner *Example*

- Counts all pairs of adjacent symbols
- The most frequent is the pair **n e**
- Merge these symbols, treating **ne** as one symbol, & add the new symbol to the vocabulary

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**

Byte-Pair-Encoding (BPE) – Token learner *Example*

vocab

30522

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_

Final vocabulary

Byte-Pair-Encoding (BPE) – Token segmenter

[coined by [Gage et al., 1994](#); adapted to the task of word segmentation by [Sennrich et al., 2016](#); see [Gallé \(2019\)](#) for more]

The token segmenter is used to tokenize a test sentence

- ❖ Just runs on the test data the merges we've learned from the training data, greedily, in the order we learned them

First we segment each test sentence word into characters

Then we apply the first merge rule

- ❖ E.g., replace every instance of **e r** in the test corpus with **er**

Then the second merge rule

- ❖ E.g., replace every instance of **er _** in the test corpus with **er_**

And so on

Vocabulary size of various models

Model	Tokenizer	Vocabulary Size
BERT base (uncased) [2018]	WordPiece ✓	30,522
BERT base (cased) [2018]	WordPiece ✓	28,996
GPT-2 [2019]	BPE ↗	50,257
Flan-T5 [2022]	SentencePiece	32,100
GPT-4 [2023]	BPE —	> 100,000
StarCoder2 [2024]	BPE ↘	49,152
Llama2 [2023]	BPE ✓	32,000

— ↗ ↘

fixed



w_1

env

done



w_2

reduced

w_3

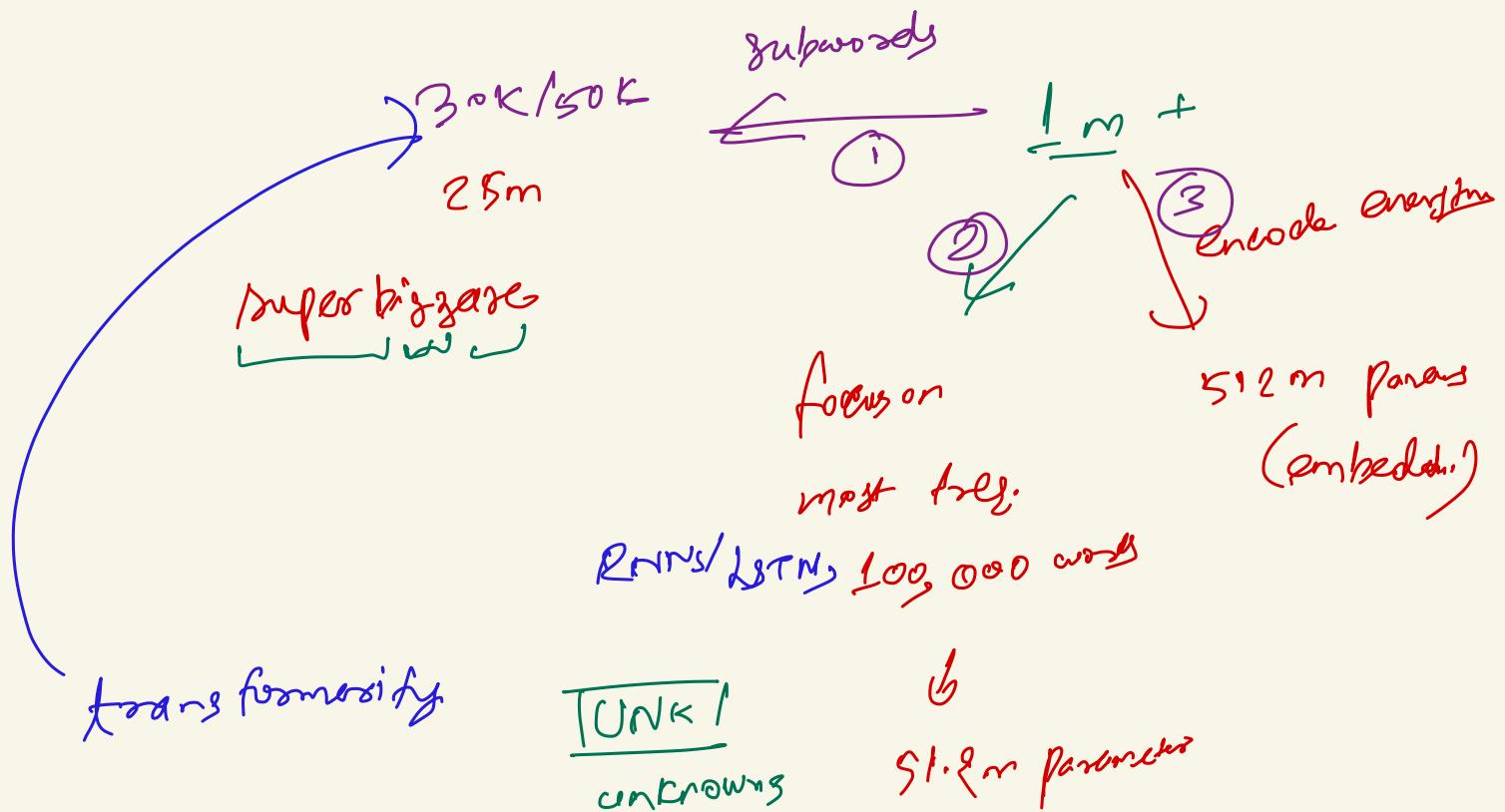
by TS

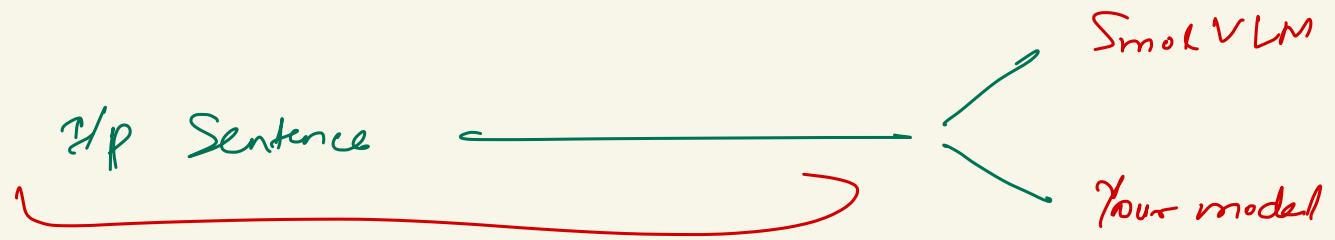
Coop(s)

S₁ Glove embedding is of 300 dimensions

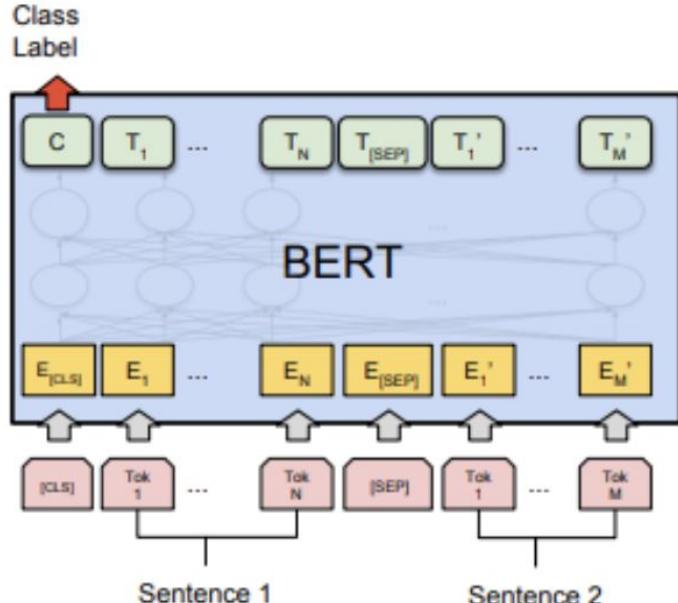
↓ BPE tokens
↑

Glove em ~~bed~~ ^{dog} _{mask} is of ~

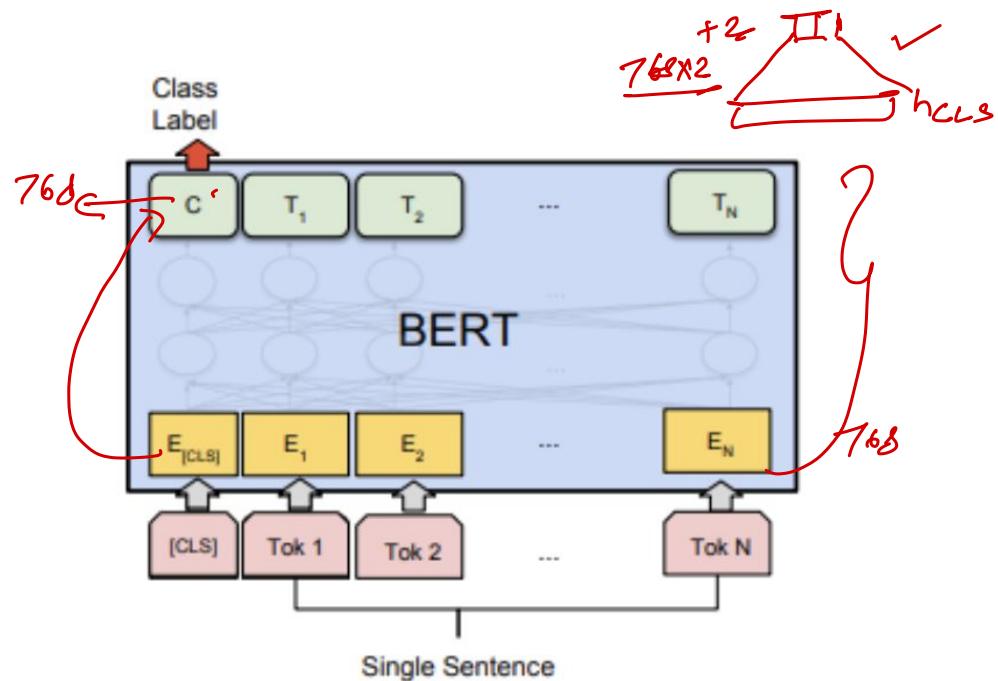




Using BERT for different tasks

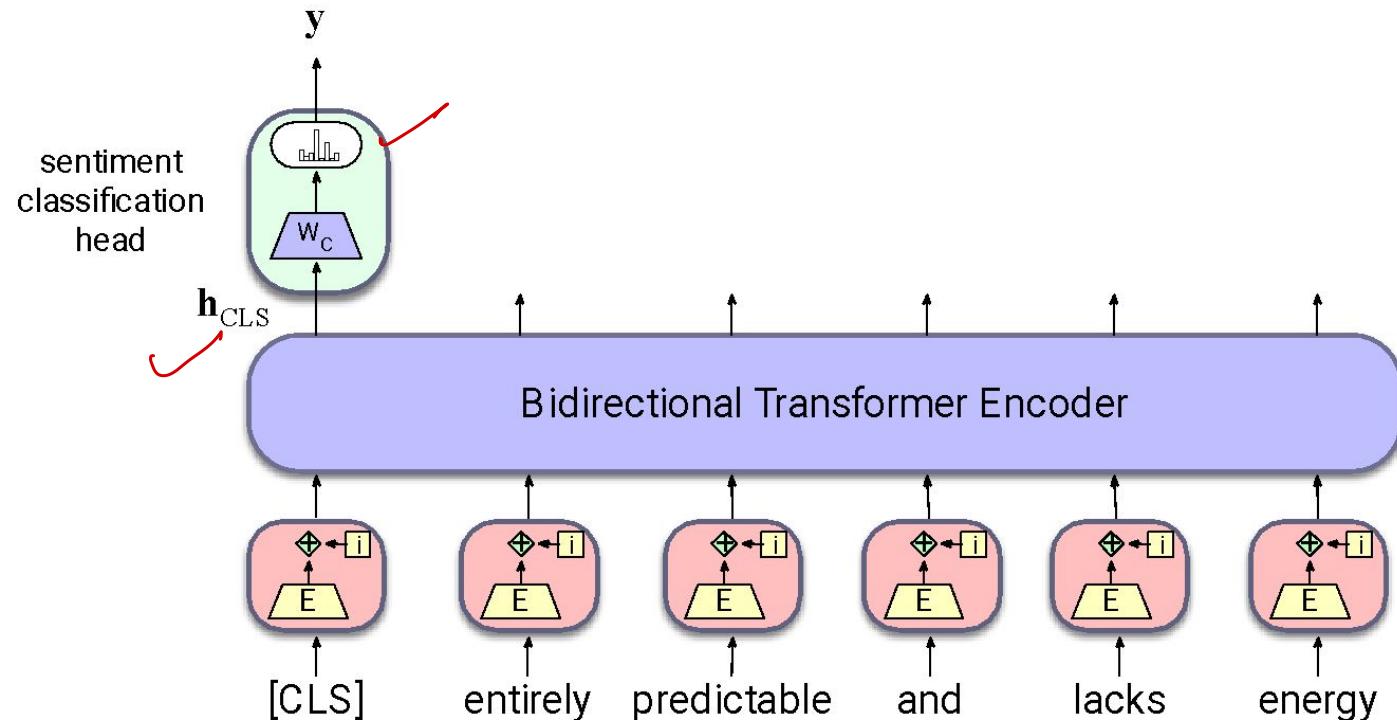


(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG

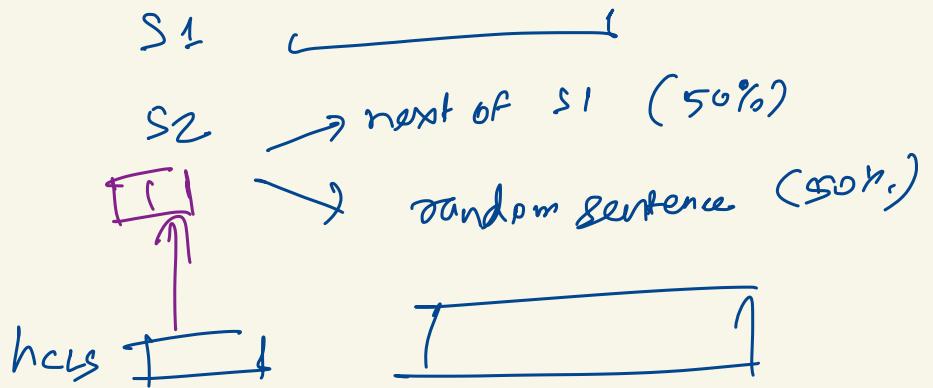


(b) Single Sentence Classification Tasks:
SST-2, CoLA

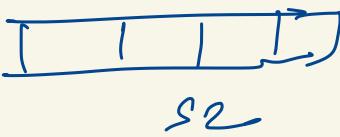
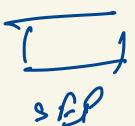
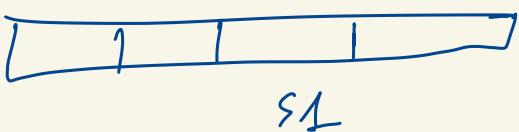
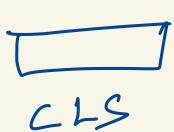
Fine Tuning for Sequence Classification



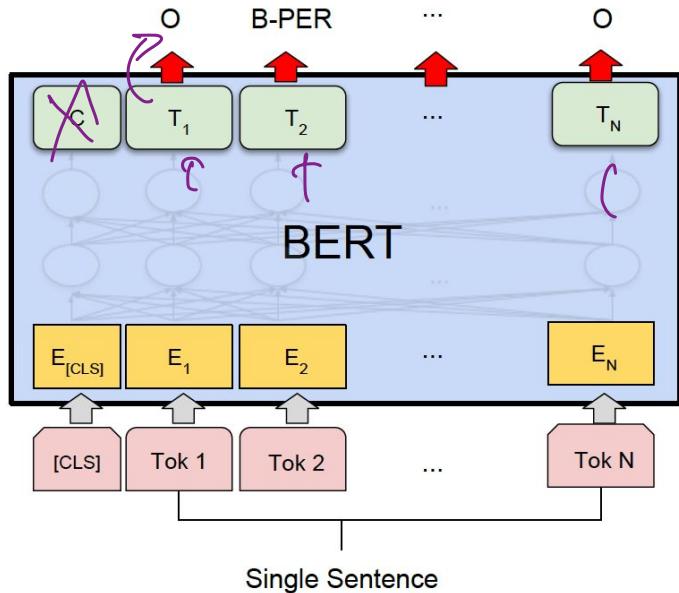
How to train CLS?



NSP (Next Sentence
Pred'g) objective



Sequence Labeling

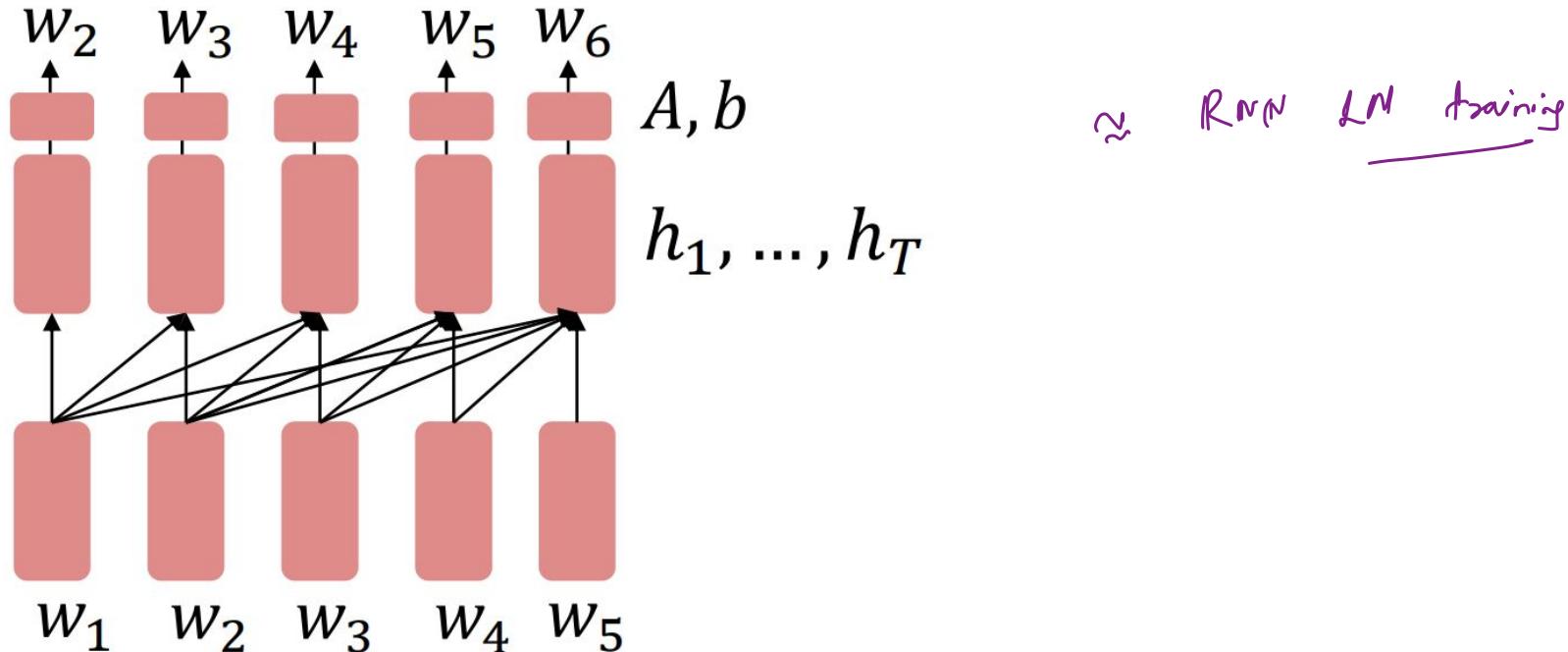


- Here, the final output vector corresponding to *each input token* is passed to a classifier that produces a softmax distribution over the possible set of tags.
- The set of weights to be learned for this additional layer is $W_K \in R^{k \times d_h}$, where k is the number of possible tags for the task.

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

Pretraining Decoders

It's natural to pretrain decoders as language models



Generative Pretrained Transformer (GPT)

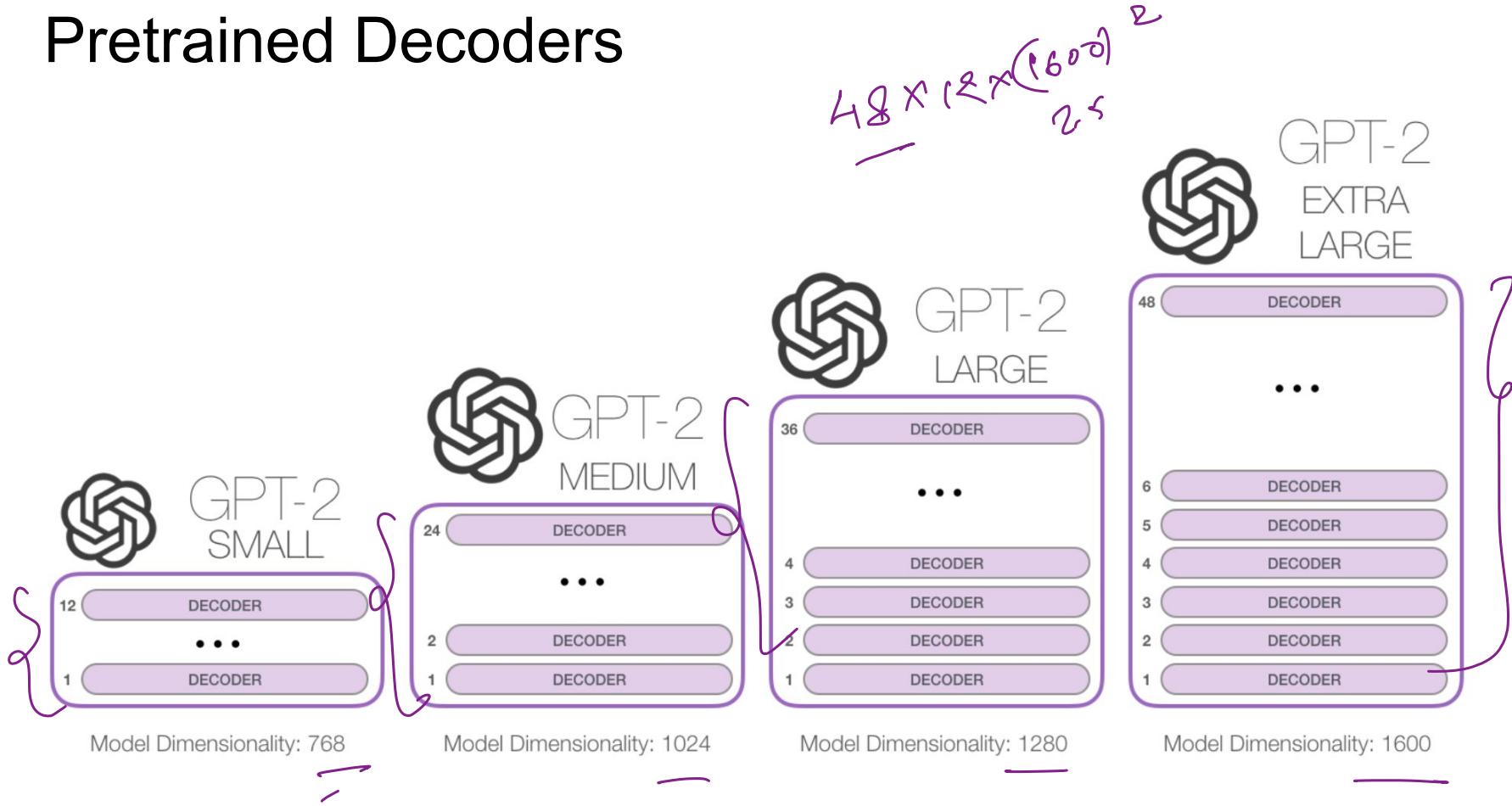
BERT

190M

- Transformer decoder with 12 layers, 117M parameters.
- 768-dimensional hidden states, 3072-dimensional feed-forward hidden layers.
30522
- Byte-pair encoding with 40,000 merges
Vocabulary larger than BERT
- Trained on BooksCorpus: over 7000 unique books.
- Contains long spans of contiguous text, for learning long-distance dependencies.

$\approx 10K \times 768 \quad \approx 7M$

Pretrained Decoders



GPT-2: Beyond Language Modeling

Basic Idea

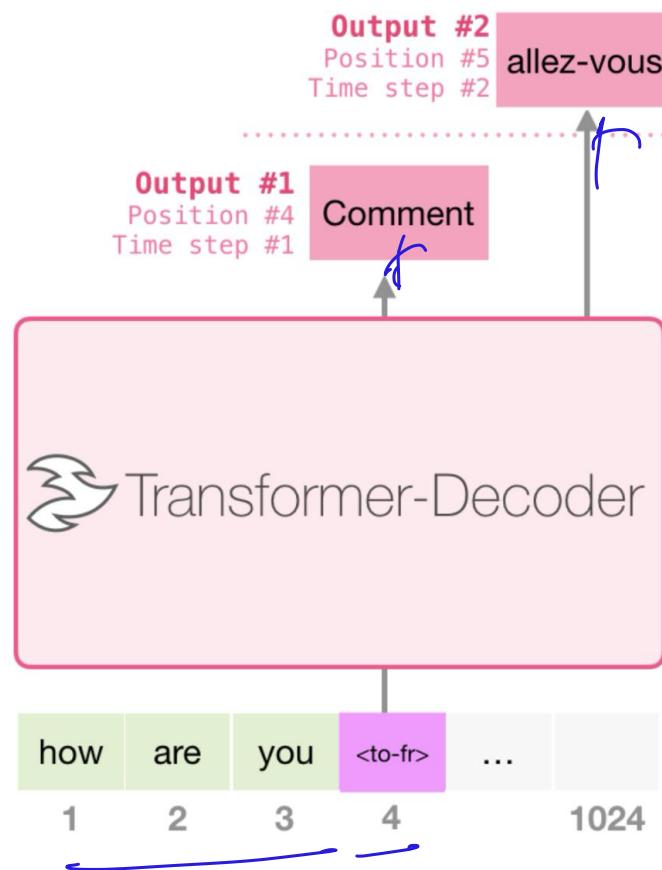
The decoder can work with a prompt!

Any NLP task can be expressed in a probabilistic framework as estimating a conditional distribution $p(\underline{output} | \underline{input})$.

GPT-2: Machine Translation

Training Dataset

I	am	a	student	<to-fr>	je	suis	étudiant
let	them	eat	cake	<to-fr>	Qu'ils	mangent	de
good	morning	<to-fr>	Bonjour				





~~Start~~

< Start >

P

Student

P

$f_{t_0-f_T}$

$\log(\tau_e)$

τ_e

P

GPT-2: Summarization

Training Dataset

Article #1 tokens	<summarize>	Article #1 Summary
Article #2 tokens	<summarize>	Article #2 Summary
Article #3 tokens	<summarize>	Article #3 Summary

