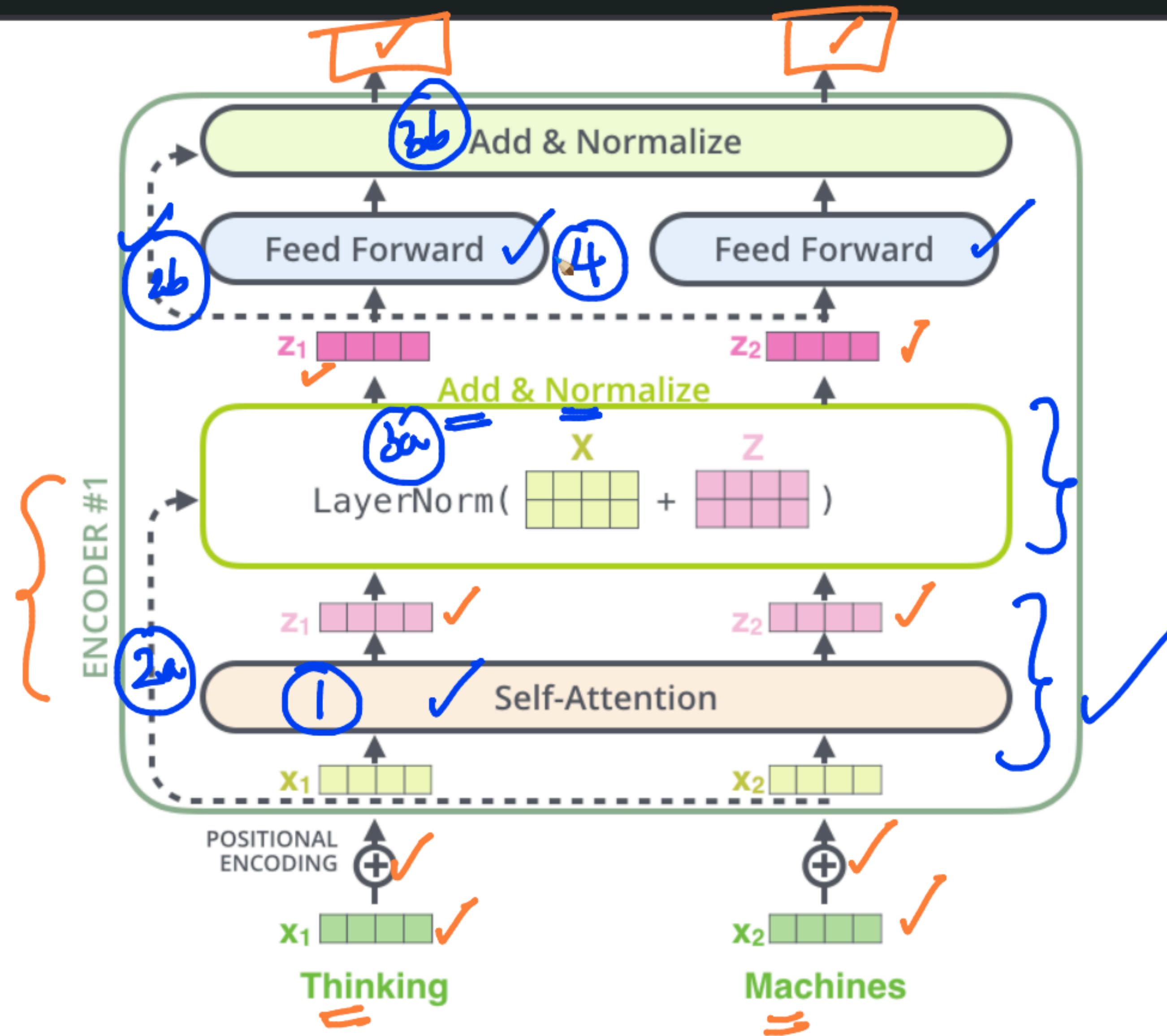
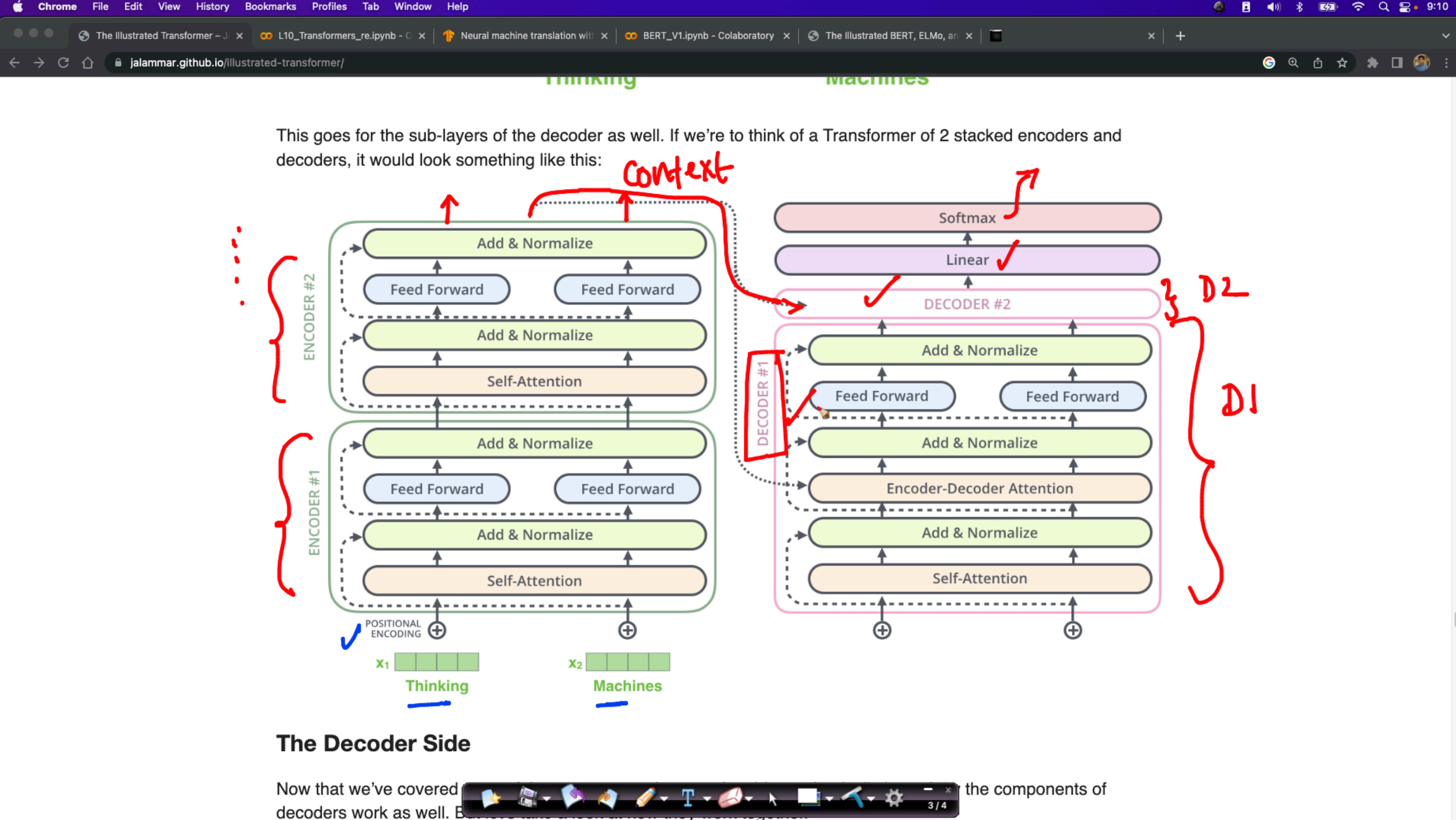


## Agenda:

1. Transformer → decoder + end-end training ✓
2. Transformer - code (TF) ✓
3. BERT - architecture ✓
4. BERT - for various NLP tasks ✓
5. BERT - code ↗

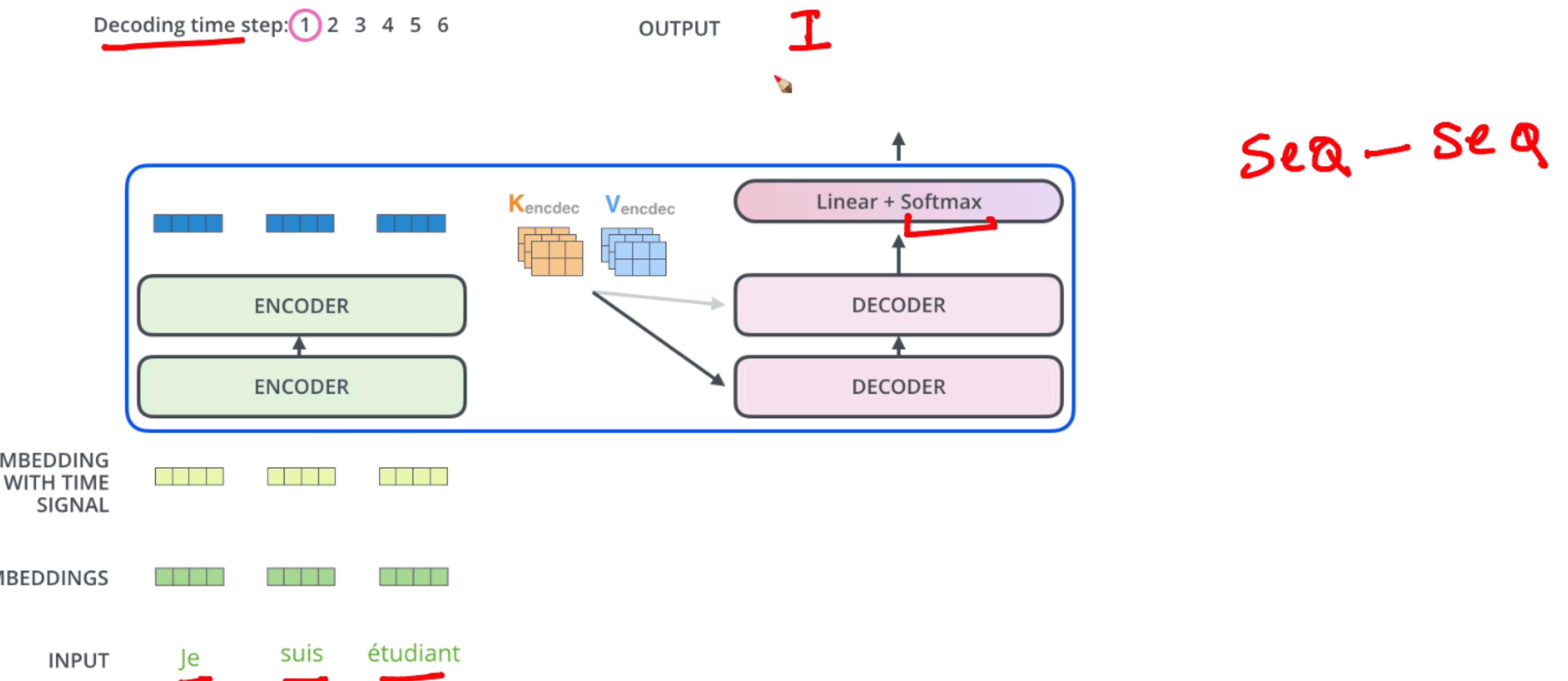


This goes for the sub-layers of the decoder as well. If we're to think of a Transformer of 2 stacked encoders and decoders, it would look something like this:



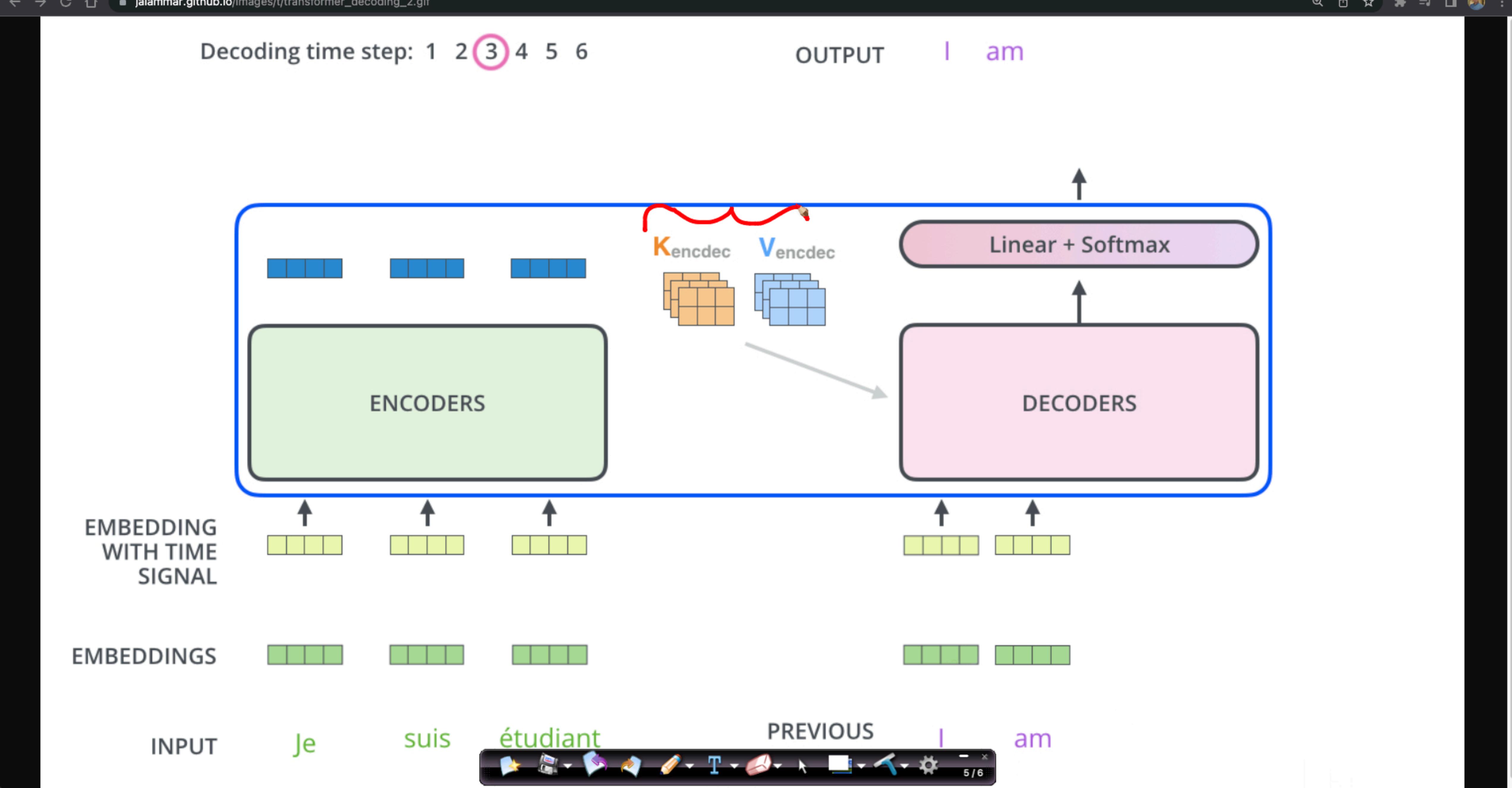
decoders work as well. But let's take a look at how they work together.

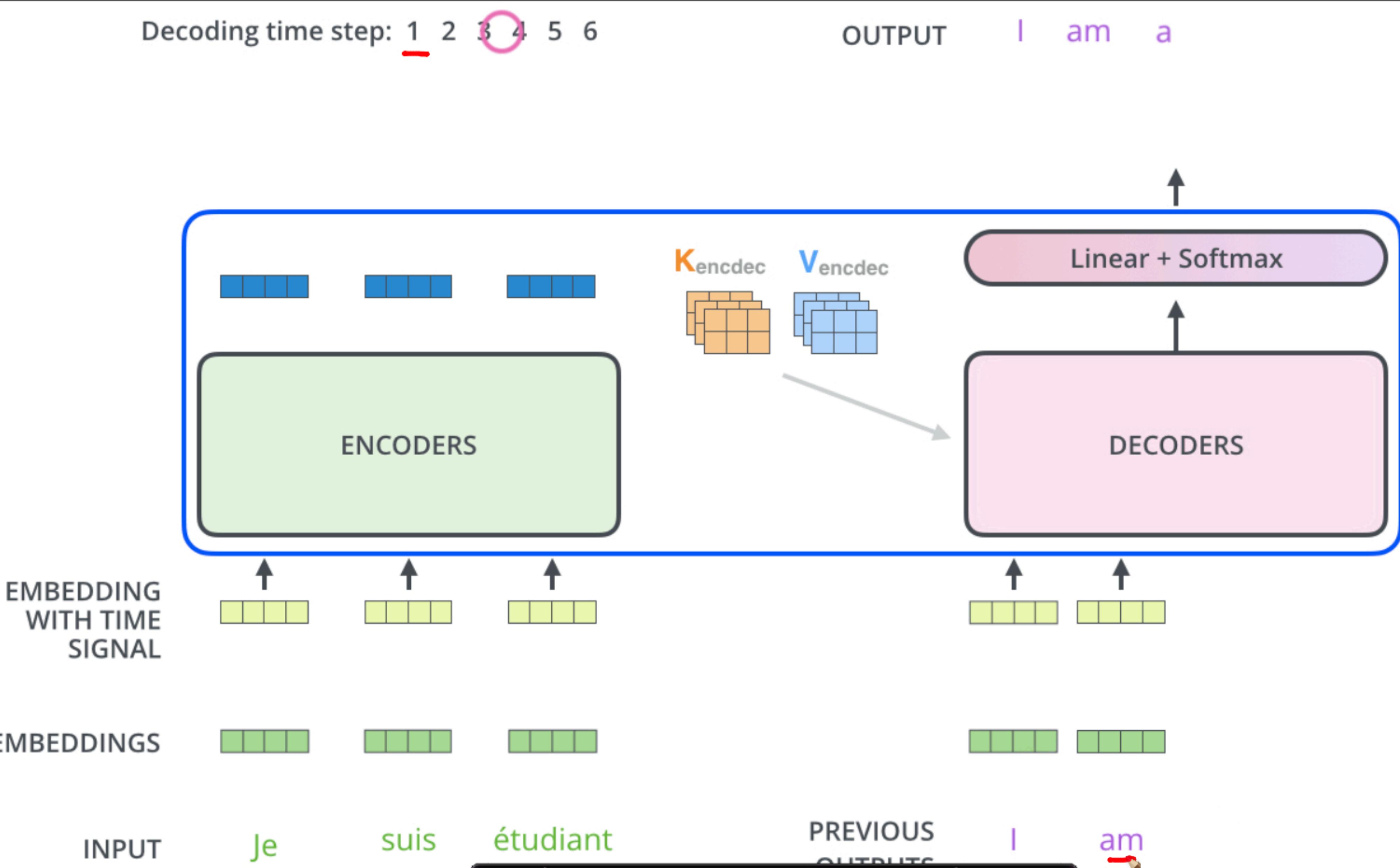
The encoder starts by processing the input sequence. The output of the top encoder is then transformed into a set of attention vectors  $K$  and  $V$ . These are to be used by each decoder in its “encoder-decoder attention” layer which helps the decoder focus on appropriate places in the input sequence:

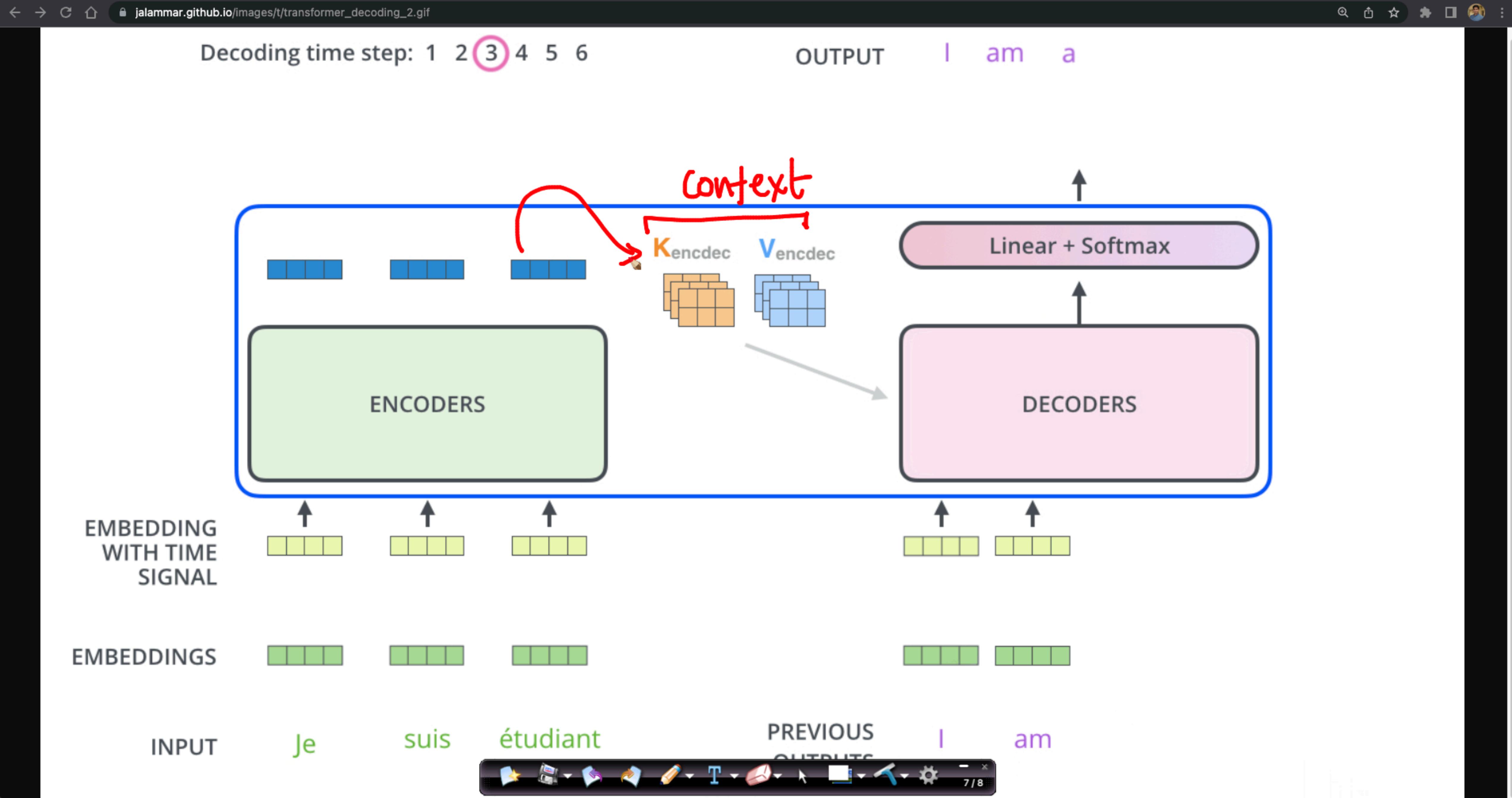


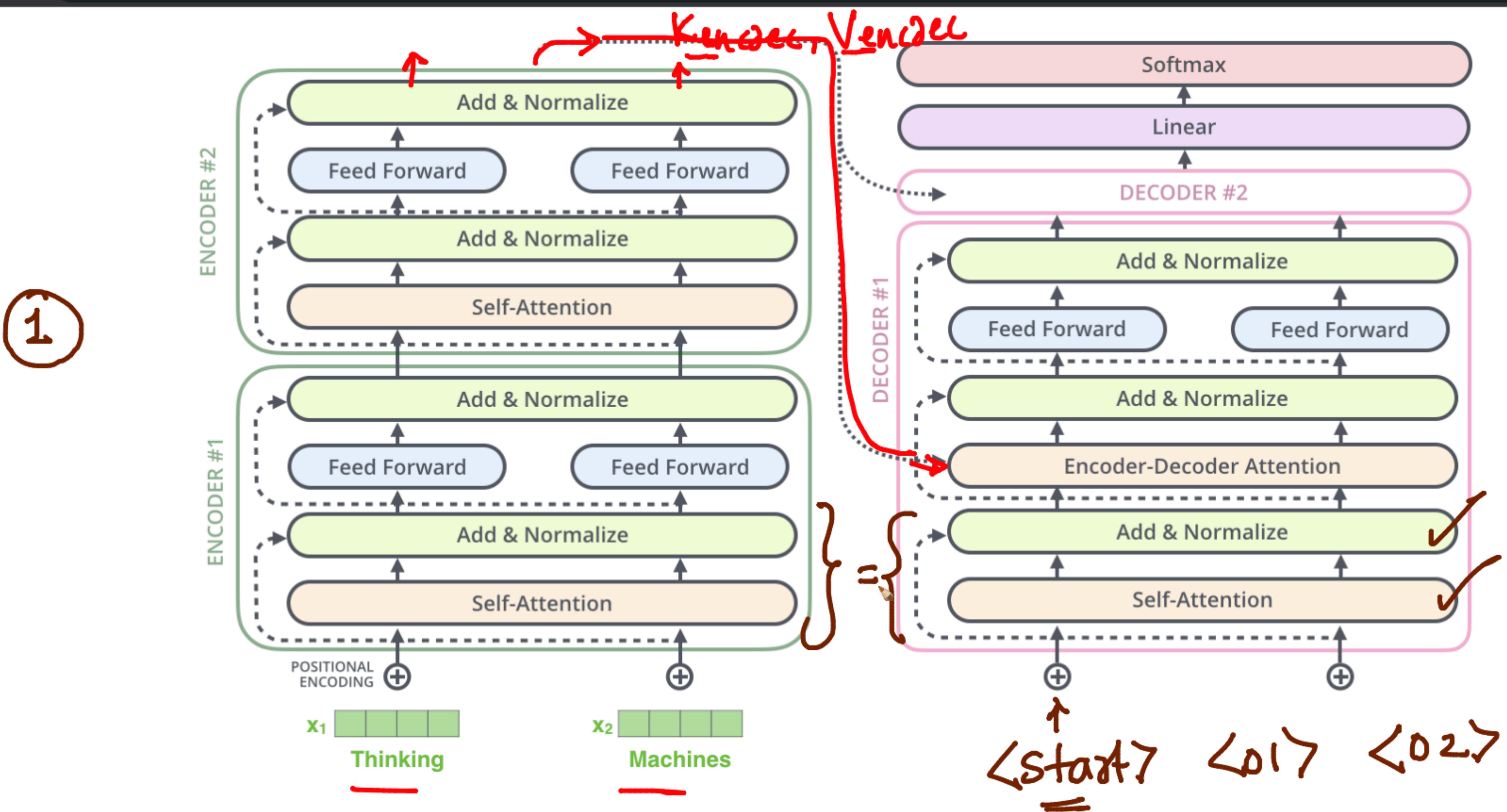
After finishing the encoding phase, we begin the decoding phase. Each step in the decoding phase outputs an element from the output sequence (the English translation sentence in this case).

The following steps repeat the process until a special symbol is reached indicating the transformer decoder has completed its output. The output of each step is fed to the bottom decoder in the next time step, and the decoders bubble up their decoding

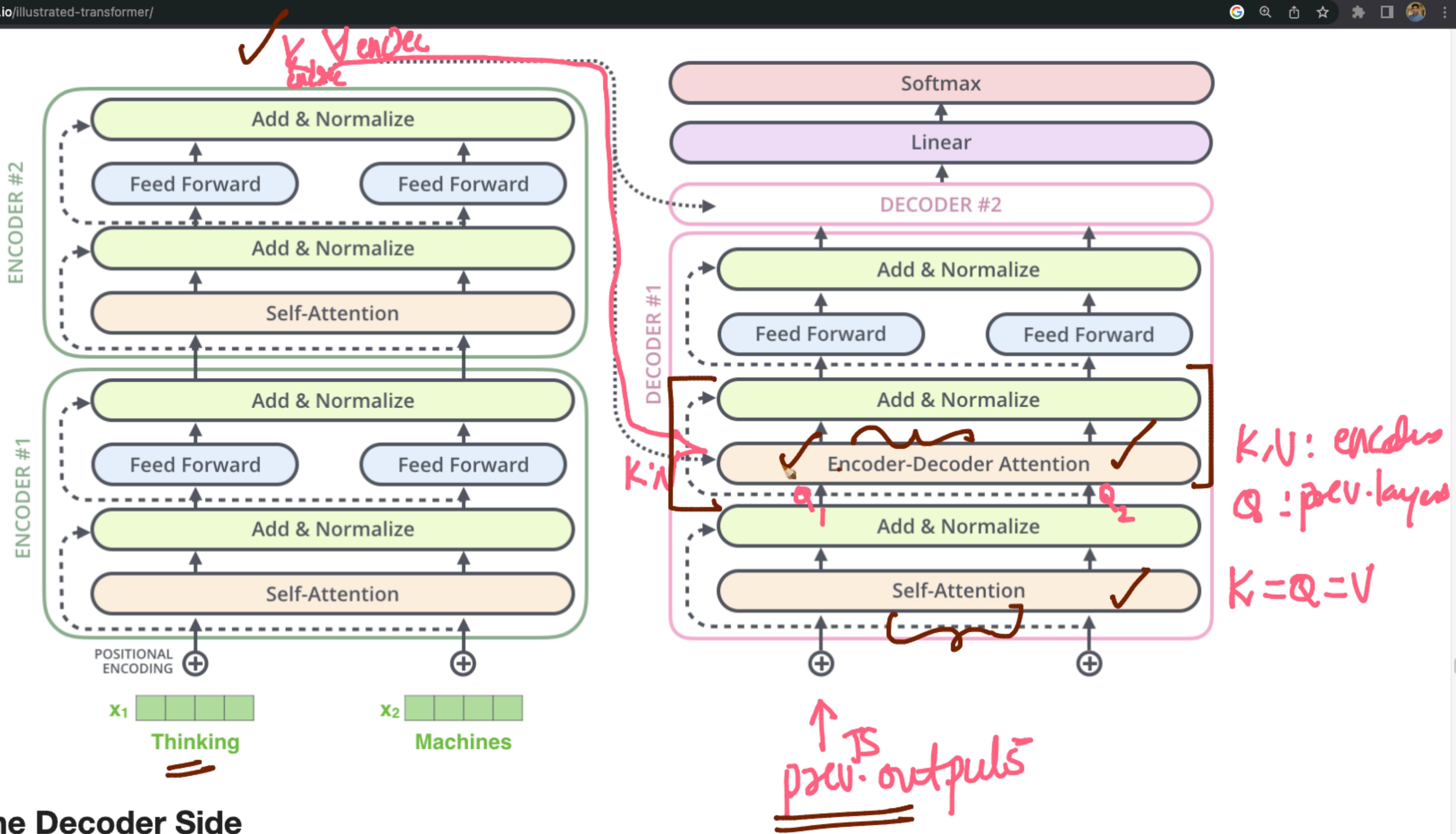




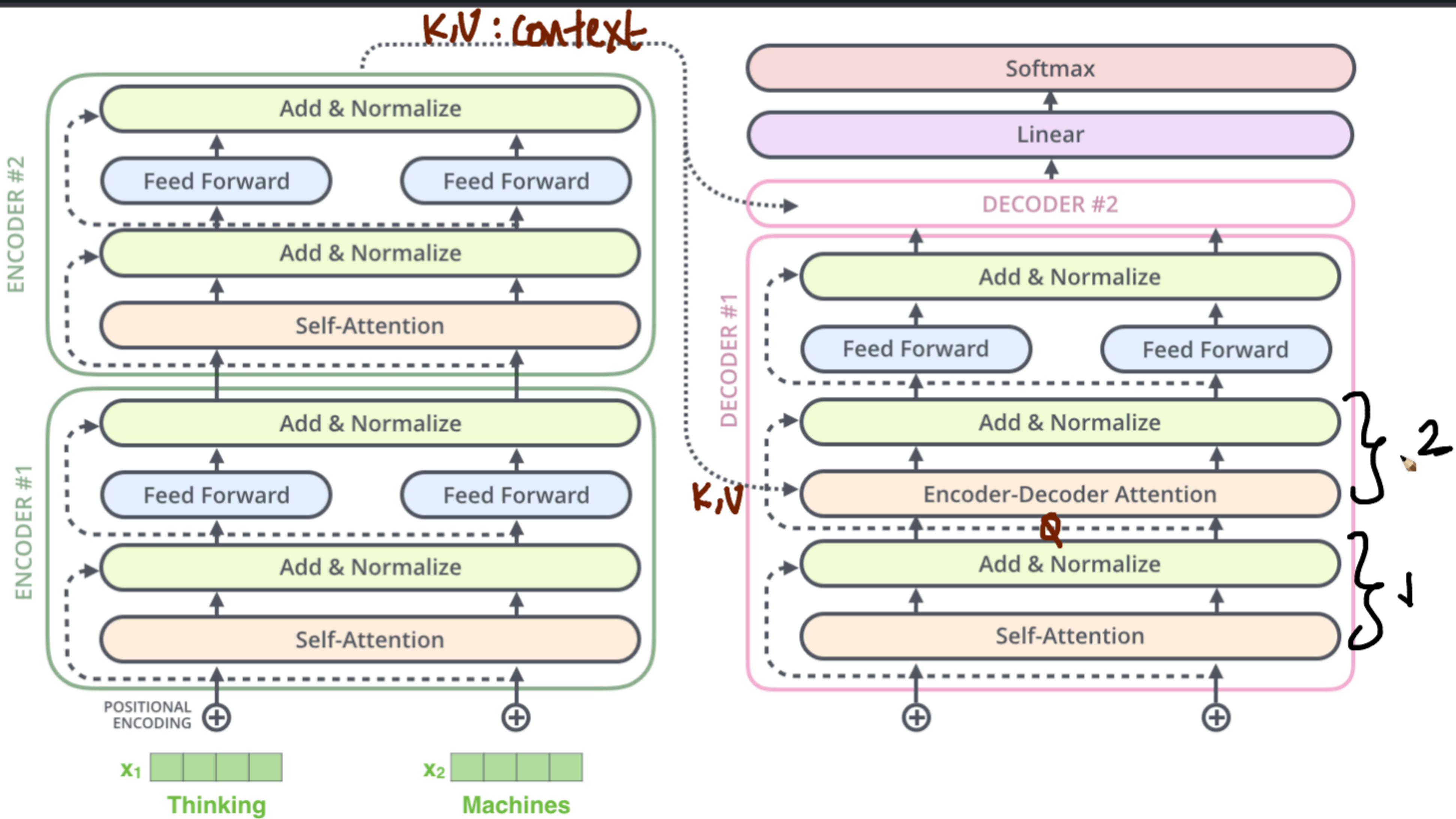




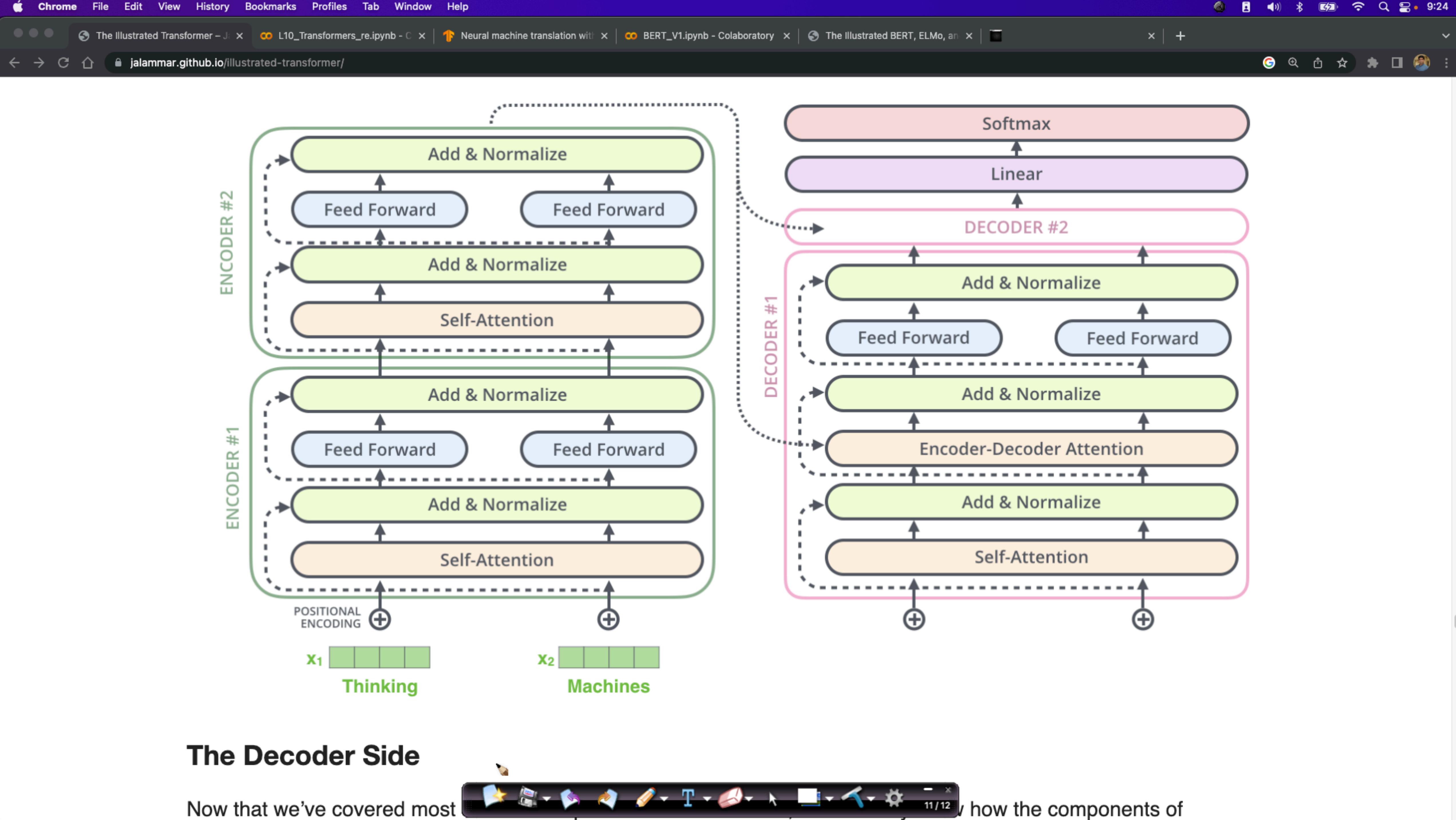
# The Decoder Side



# The Decoder Side



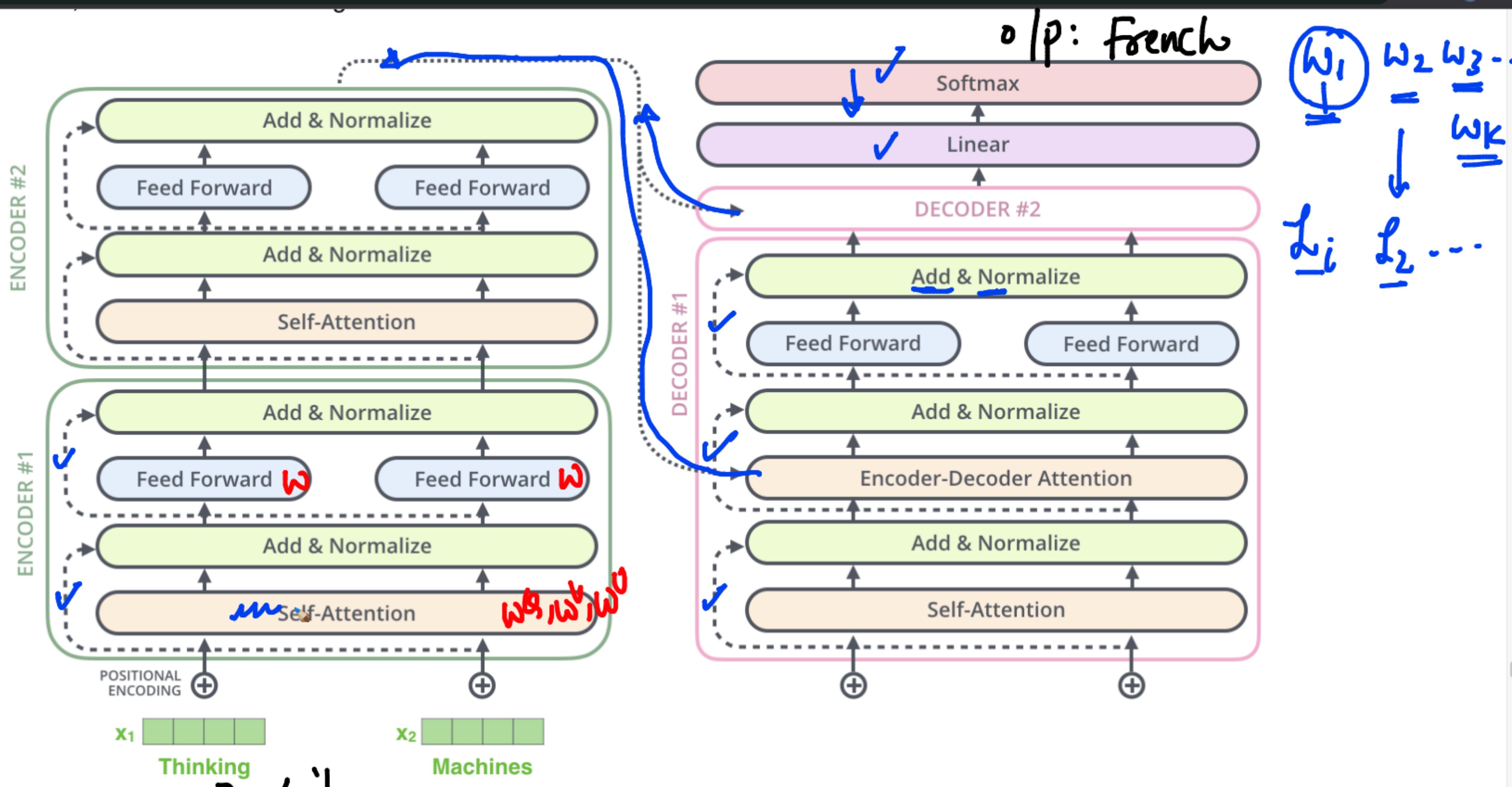
# The Decoder Side



## The Decoder Side

Now that we've covered most



~~end-end~~**The Decoder Side**Now that we've covered most of the components on the encoder side, let's call them by their common name: **transformer blocks**.

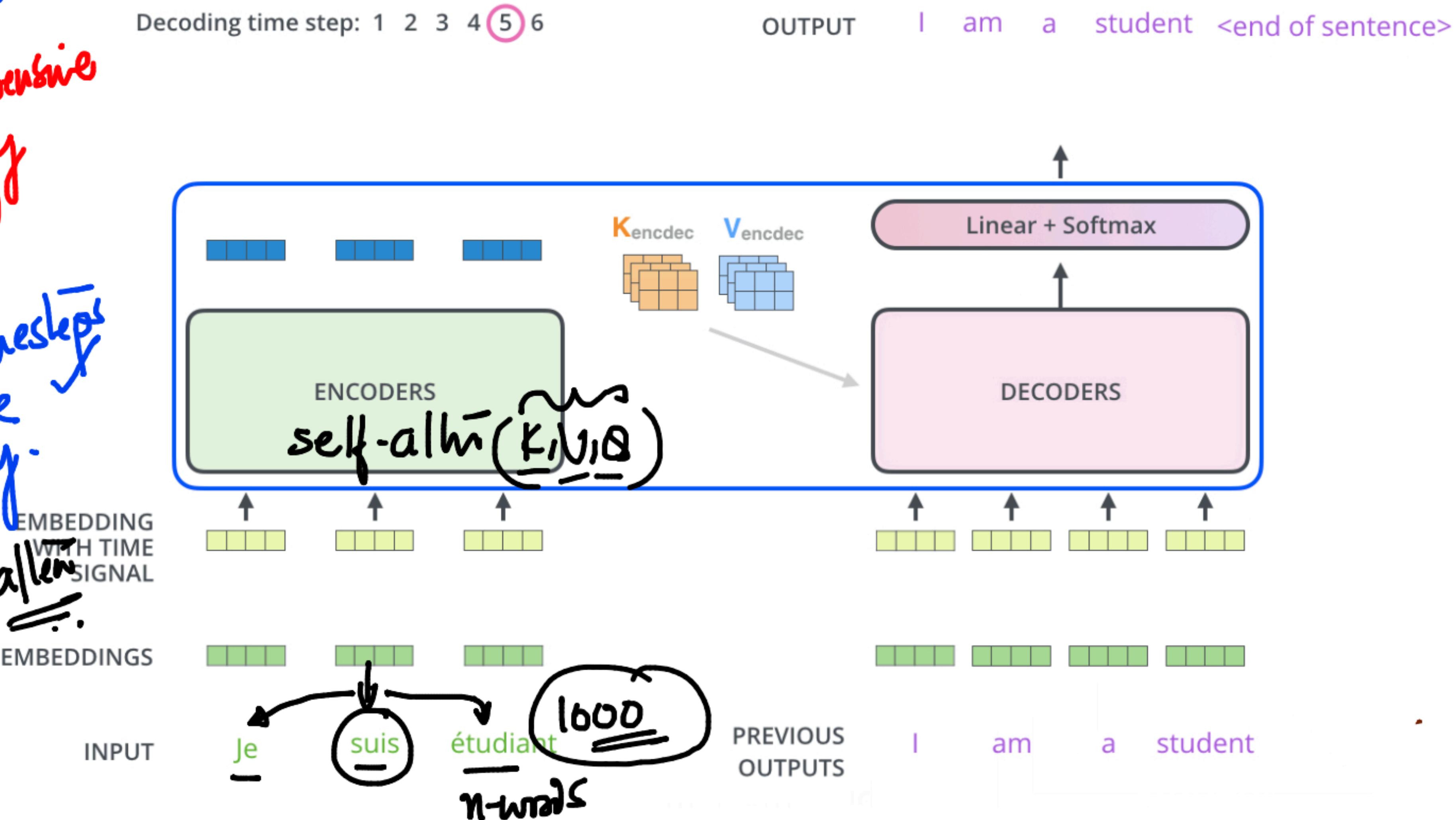
add positional encoding to those decoder inputs to indicate the position of each word.

*dauback(s):*

→ COMP. EXPENSIVE

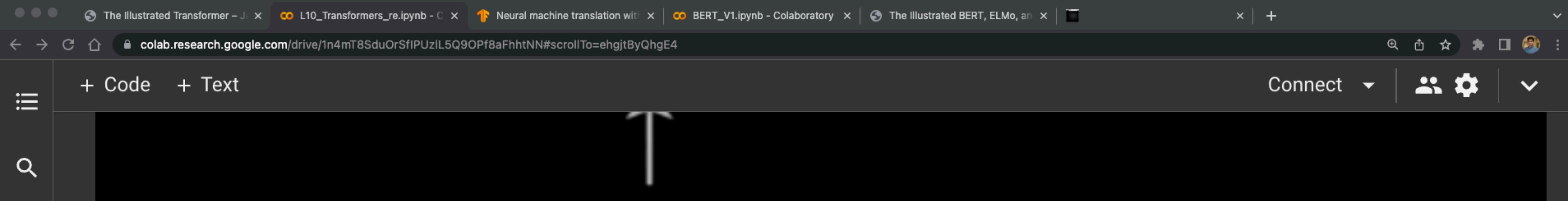
- dec. timestep  
could be many.

WITH TIME  
SIGNAL  
EMBEDDINGS



The self attention layers in the decoder operate in a slightly different way than the one in the encoder:

In the decoder, the self-attention layer takes the output sequence. This is



{x}

## □ ▾ Can you use this architecture for translation problem?

Lets see Time Complexity between RNN vs Attention

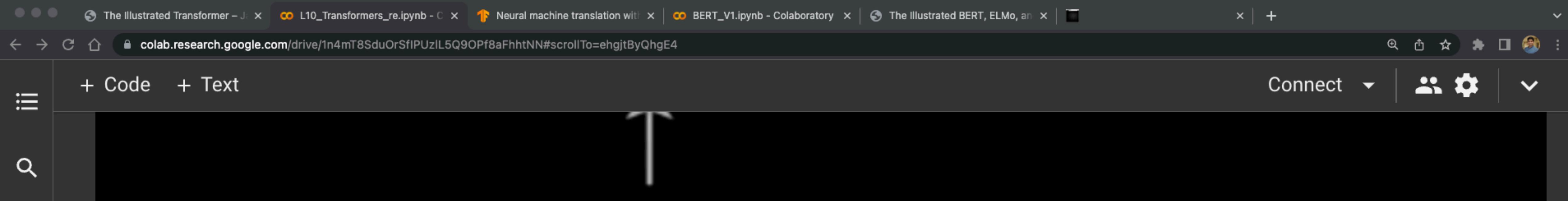
Architecture	Time Complexity per Layer	Sequential operations
RNN/LSTM	$O(n \cdot d^2)$	$O(n)$
Attention	$O(n^2 \cdot d)$	$O(1)$

- n: is the length of the sequence.
- d: is the dimension of the word representation (512, 1024 in generally)

In terms of computational complexity,

- when  $n < d$ , the attention layers are faster than the recursive layers

## Quality comparison



{x}

## Can you use this architecture for translation problem?

Lets see Time Complexity between RNN vs Attention

Architecture	Time Complexity per Layer	Sequential operations
RNN/LSTM	$O(n \cdot d^2)$	$O(n)$
Attention	$O(n^2 \cdot d)$	$O(1)$

- n: is the length of the sequence.
- d: is the dimension of the word representation (512, 1024 in generally)

In terms of computational complexity,

- when  $n < d$ , the attention layers are faster than the recursive layers

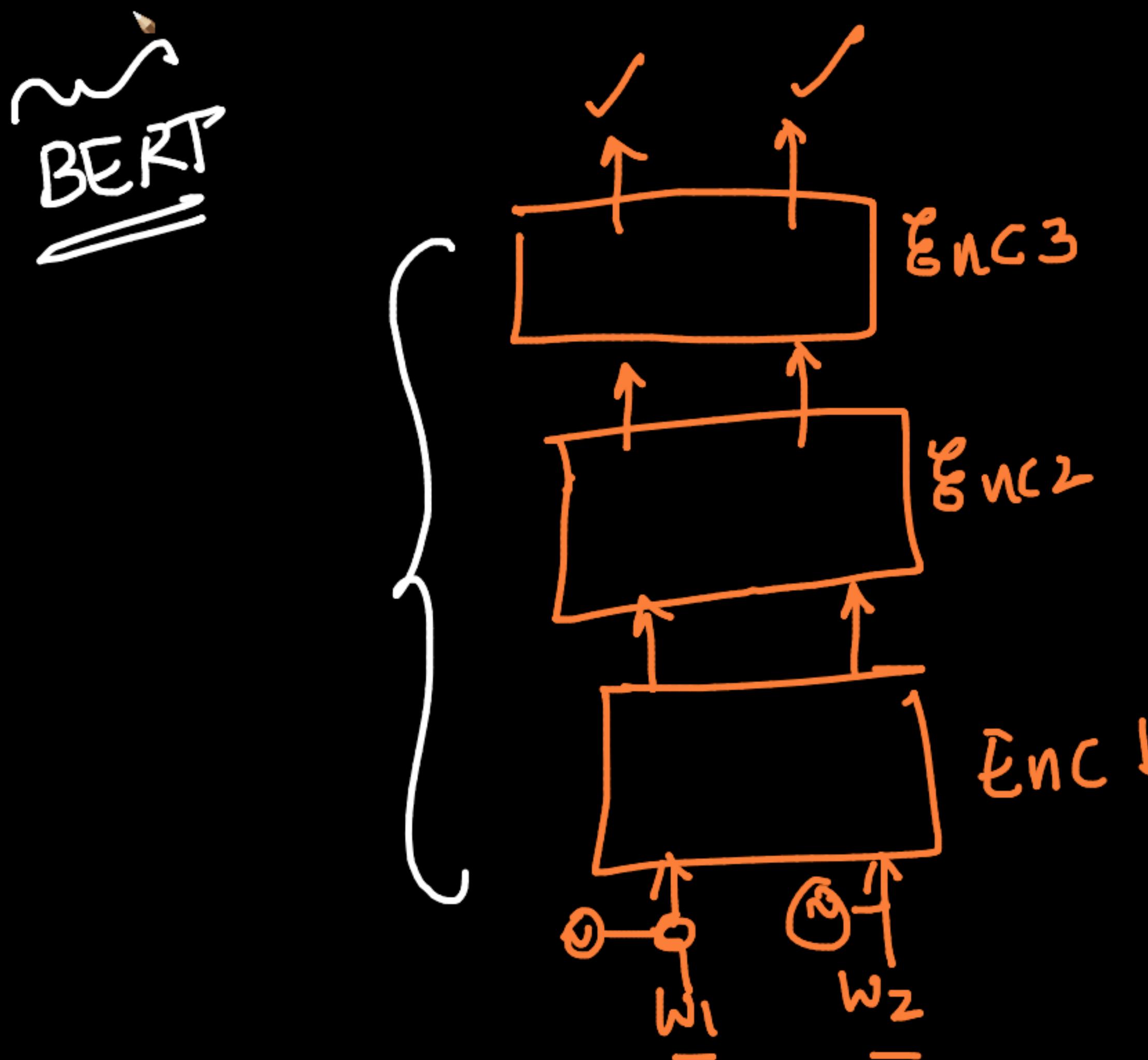
## Quality comparison

Transformer → ~~decoding per time step~~

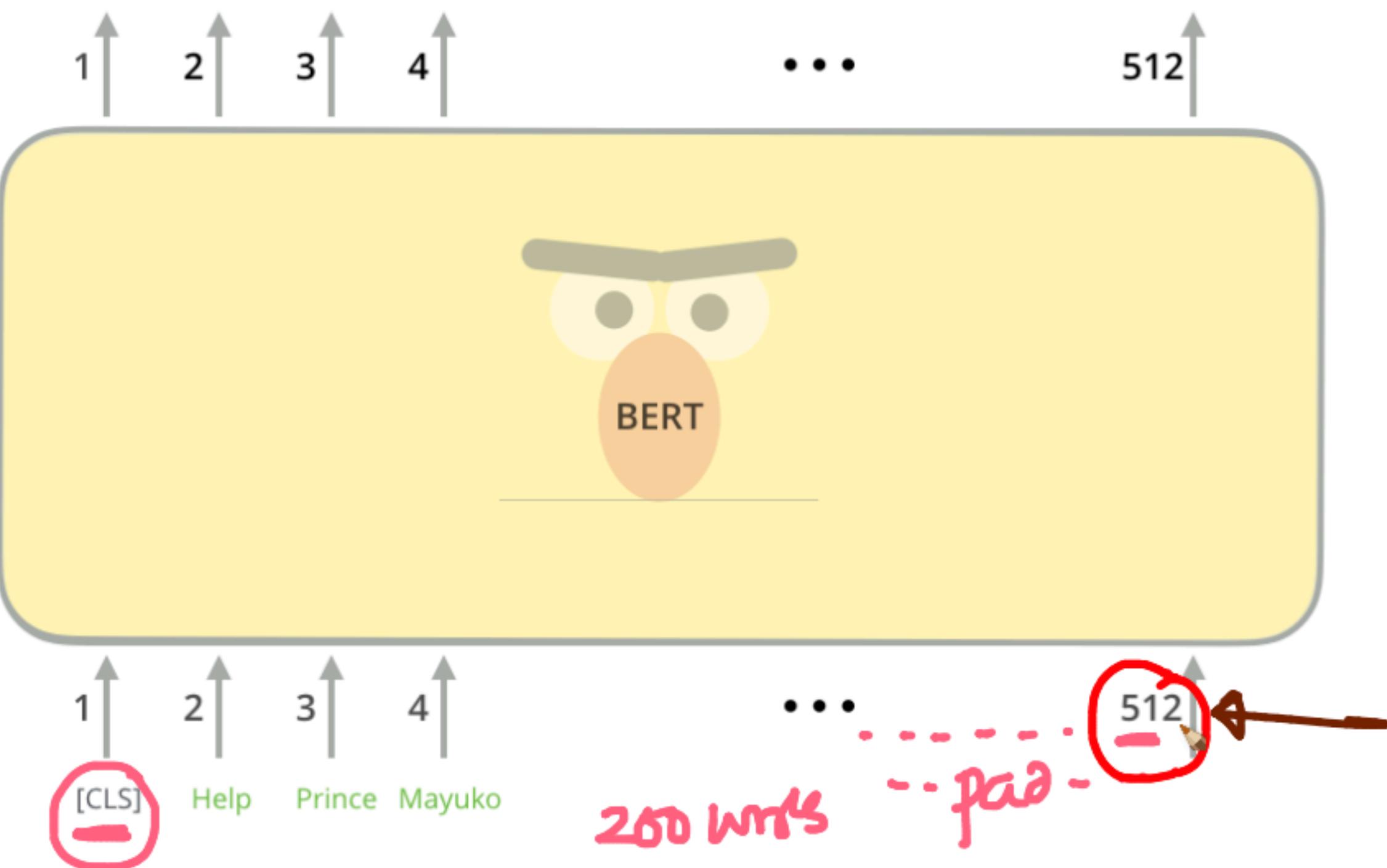


NLP:

- ✓ ① Classfn (Multi-class)
- ✓ ② Seq2Seq



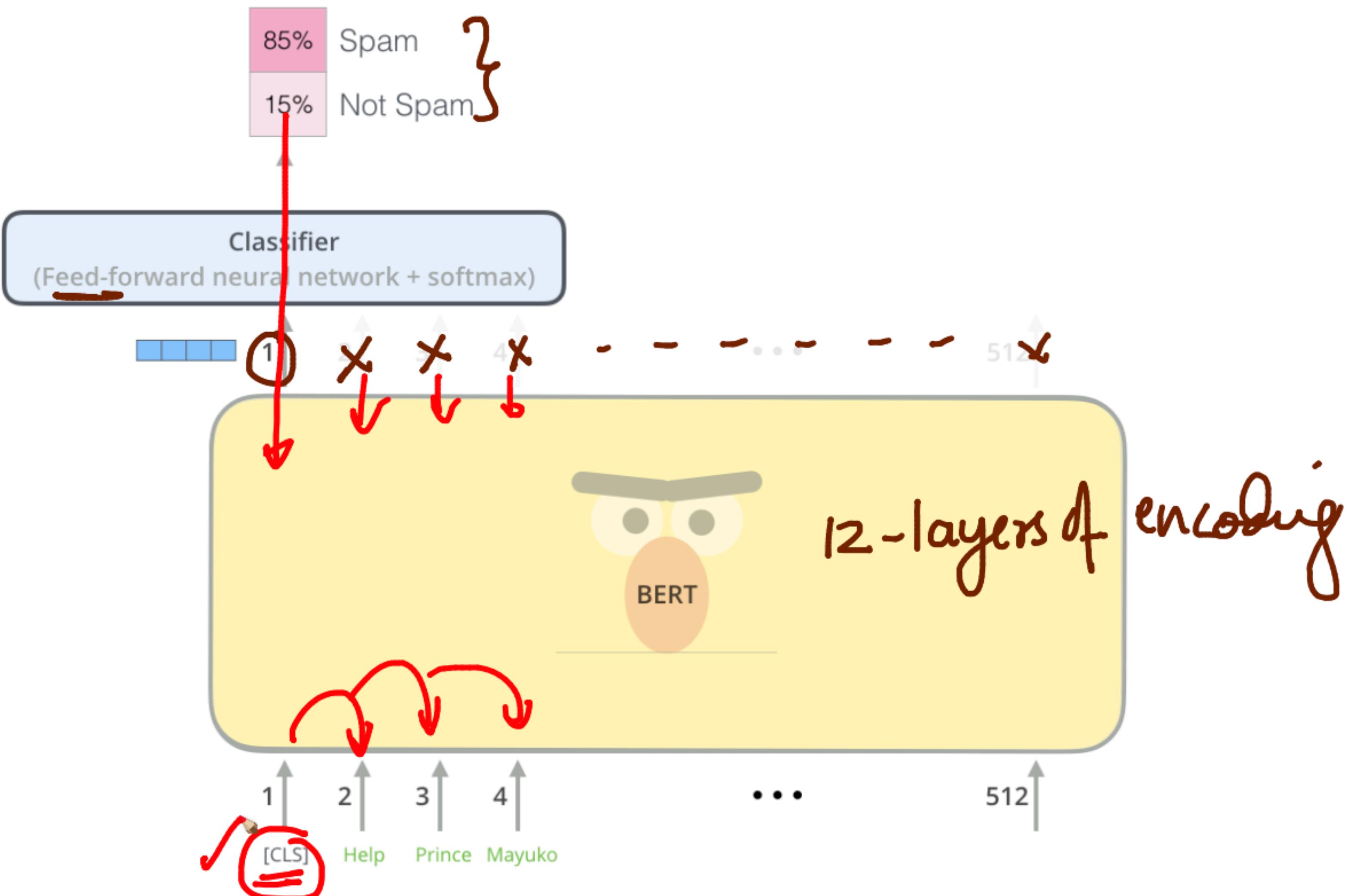
## Model Inputs



The first input token is supplied with a special [CLS] token for reasons that will become apparent later on. CLS here stands for Classification.

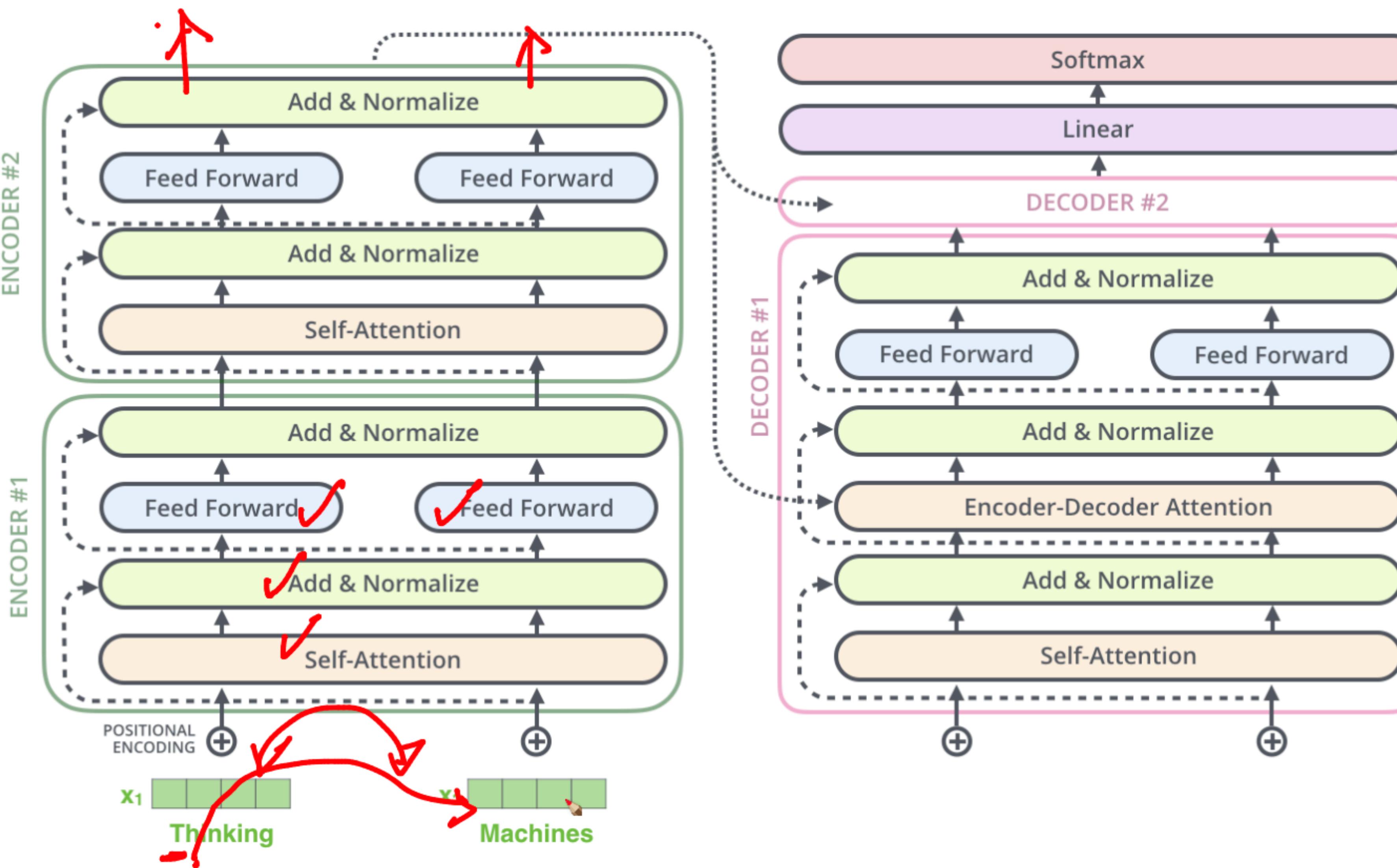
Just like the vanilla encoder of the transformer, BERT takes a sequence of words as input which keep flowing up the stack. Each layer applies self-attention, and passes its results through a feed-forward network, and then hands it off to the next encoder.

That vector can now be used as the input for a classifier of our choosing. The paper achieves great results by just using a single-layer neural network as the classifier.

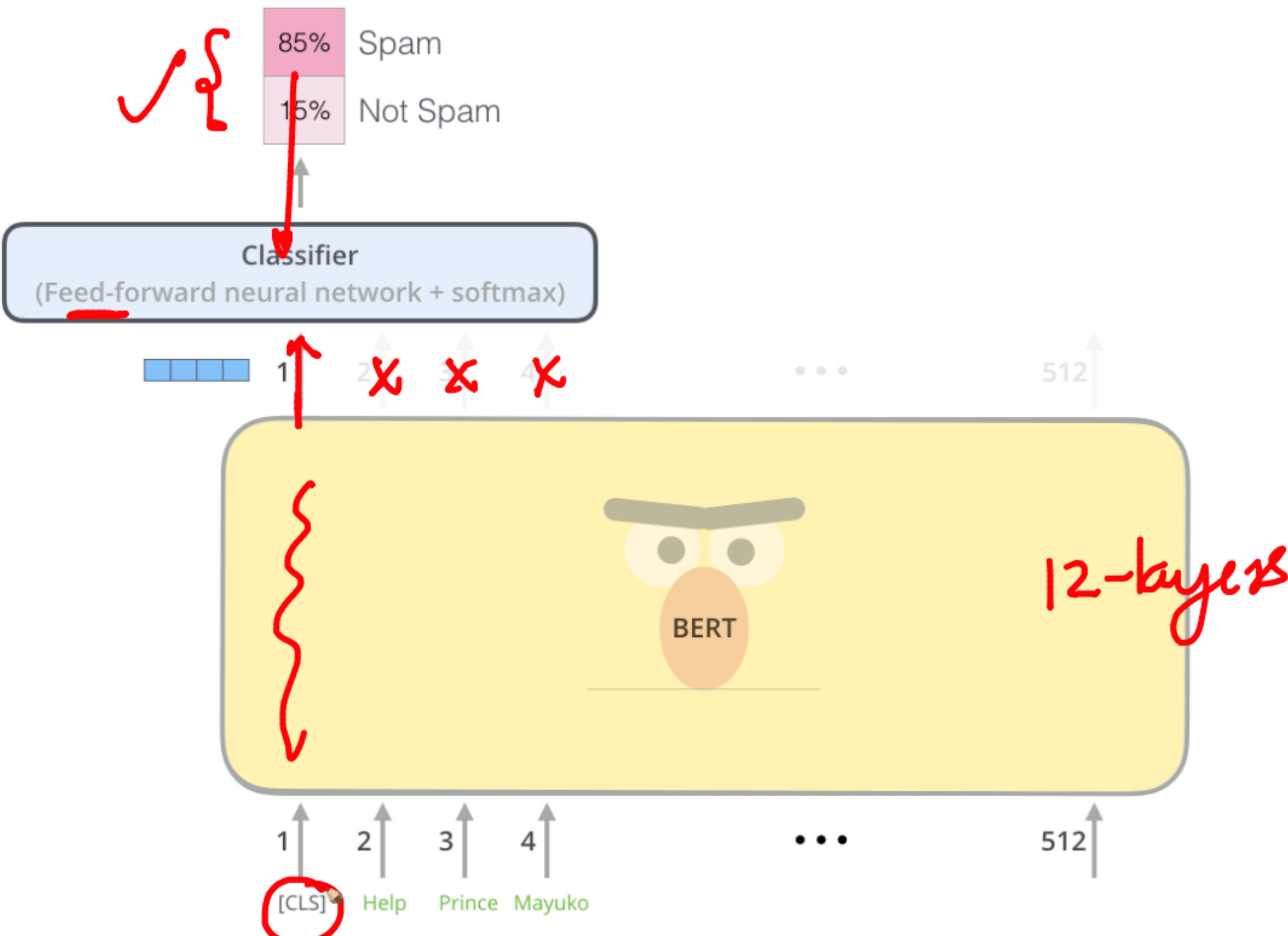


If you have more labels (for example, "spam", "not spam", "social", and "promotion"), you just tweak the classifier you work on to have more output classes and then pass through softmax.

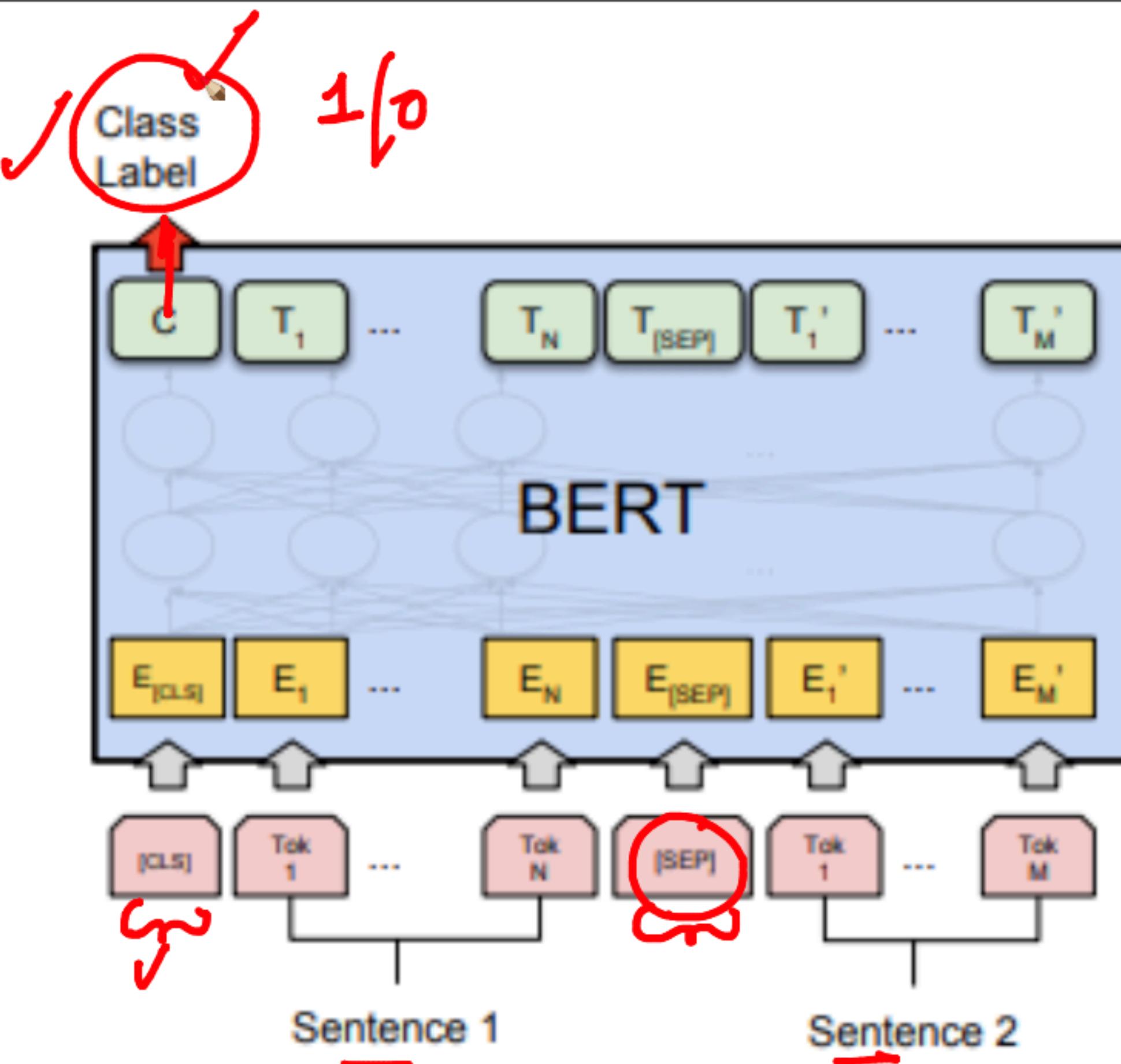
This goes for the sub-layers of the decoder as well. If we're to think of a Transformer of 2 stacked encoders and decoders, it would look something like this:



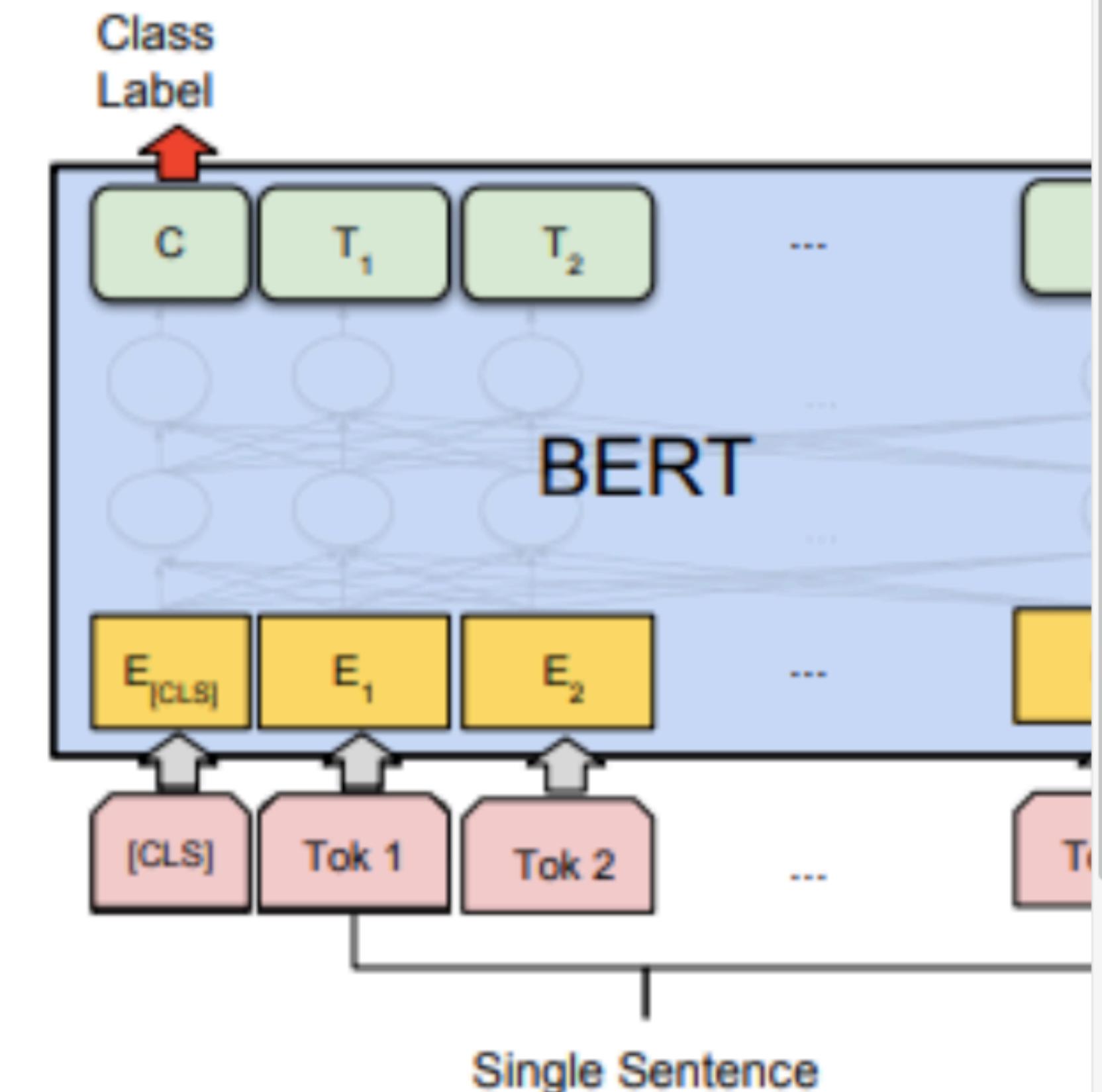
That vector can now be used as the input for a classifier of our choosing. The paper achieves great results by just using a single-layer neural network as the classifier.



If you have more labels (for example "spam", "not spam", "social", and "promotion") you just tweak the crossover now we're going to go through softmax.



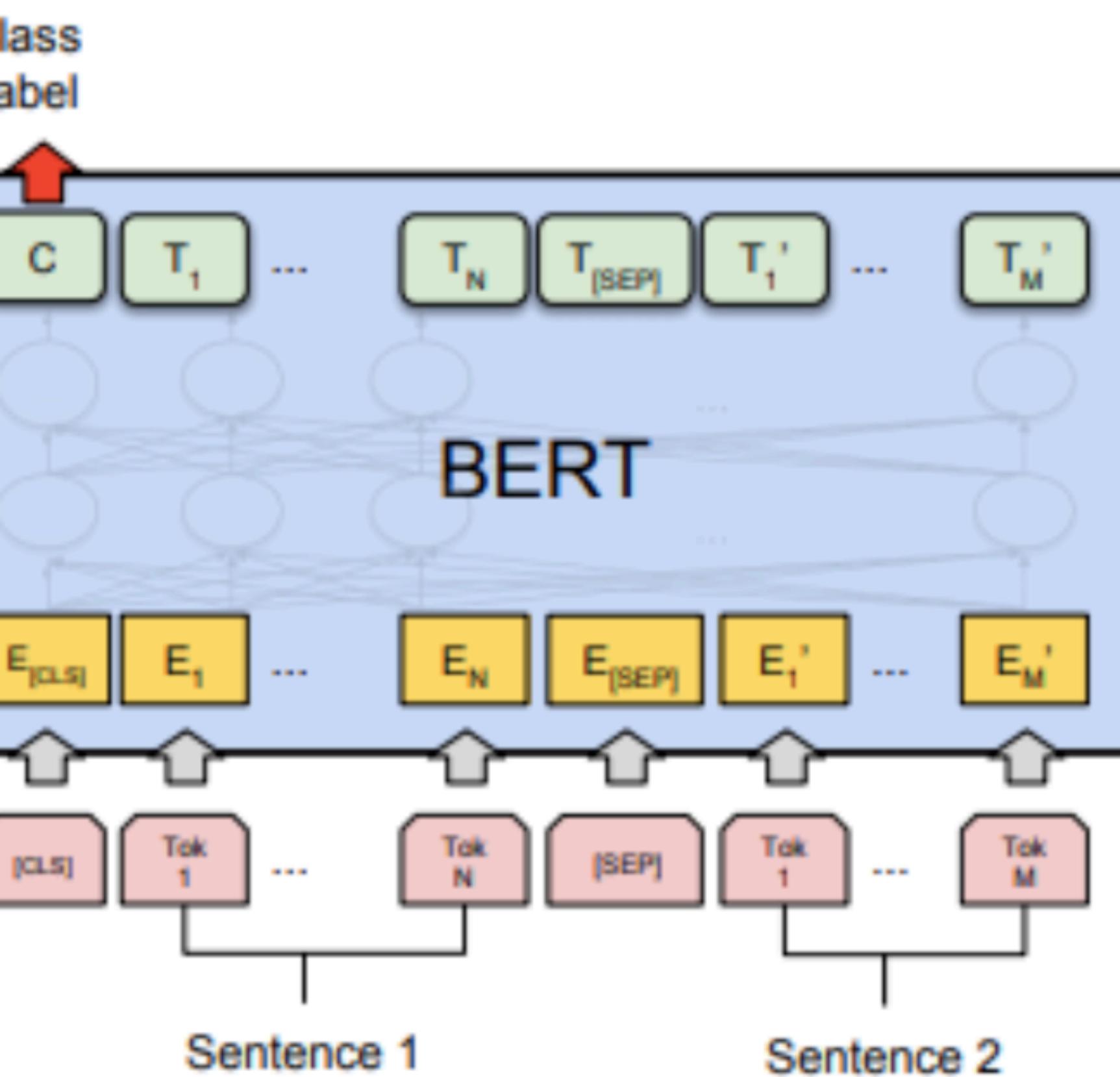
(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG



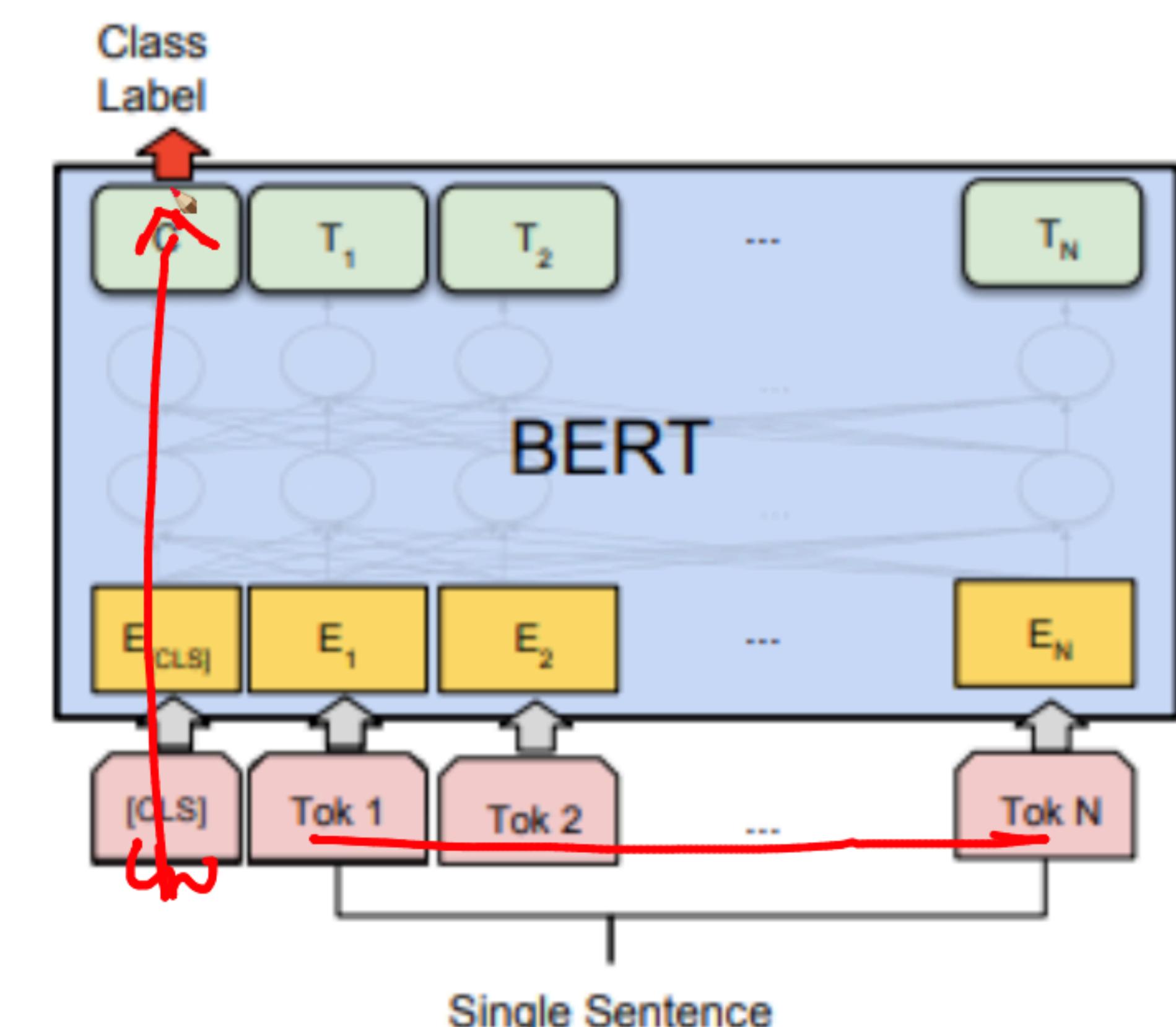
(b) Single Sentence Classification Tasks  
SST-2, CoLA

Start/End Span





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



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

Start/End Span

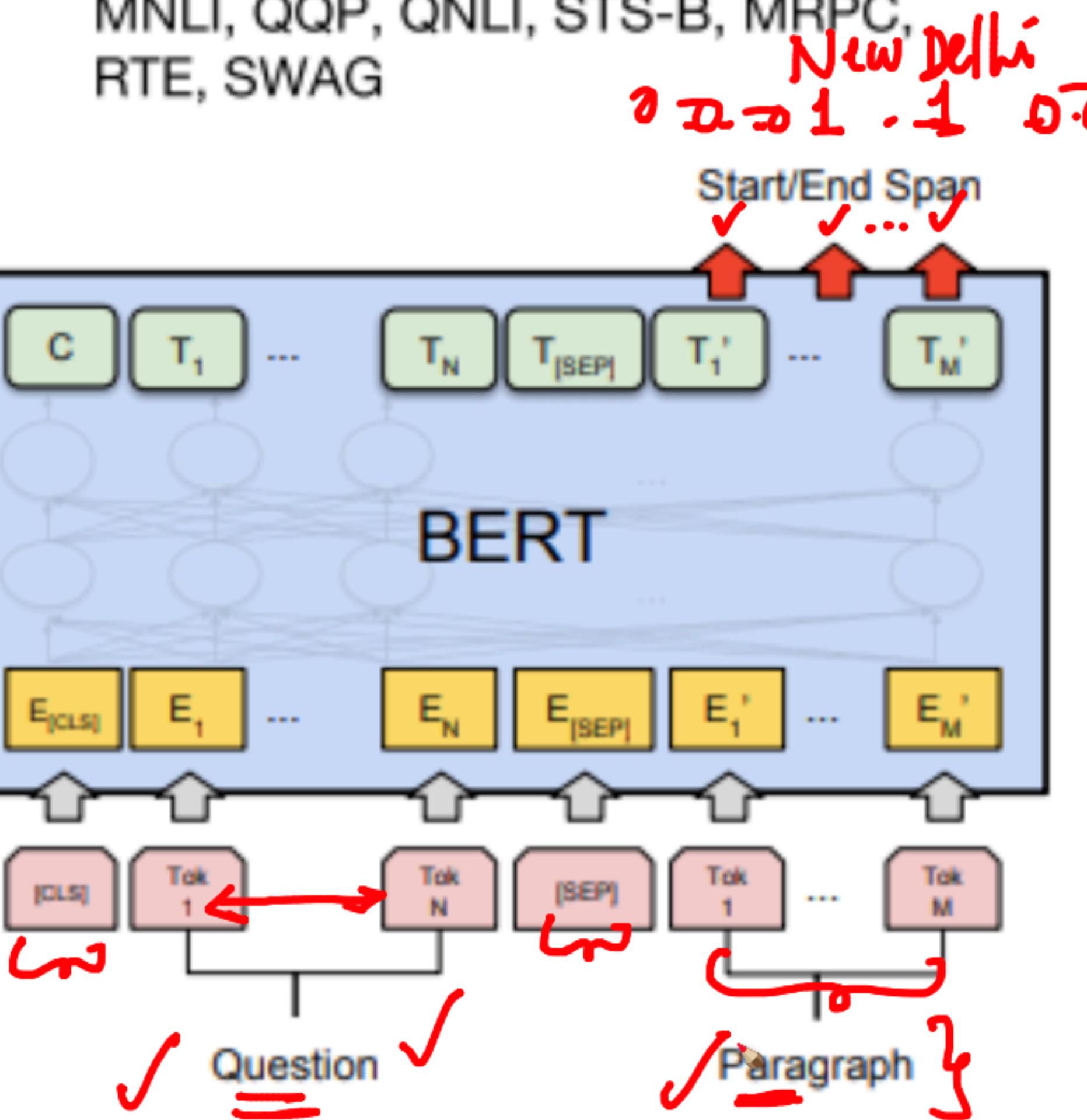
O

B-PER

...

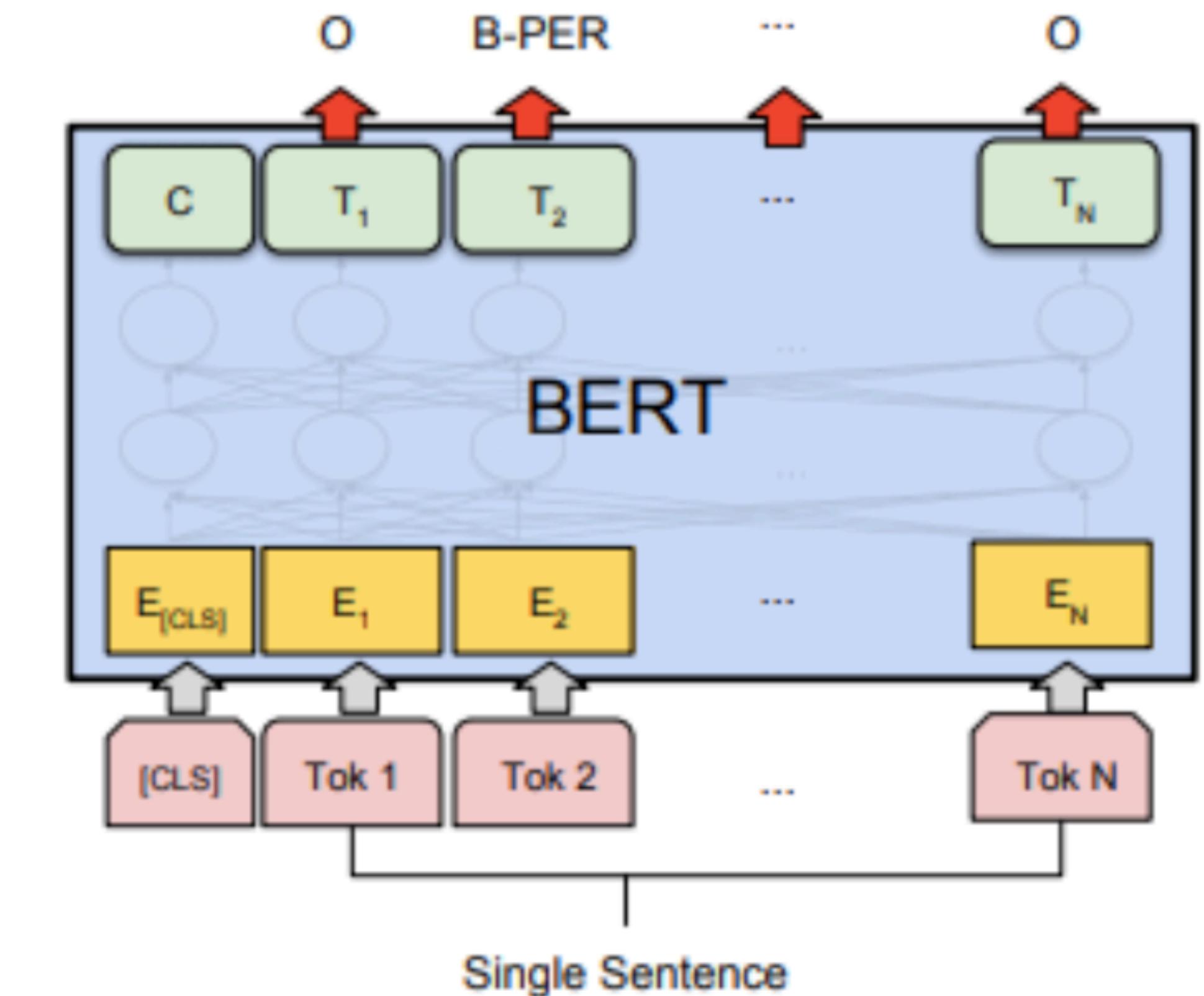
O

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



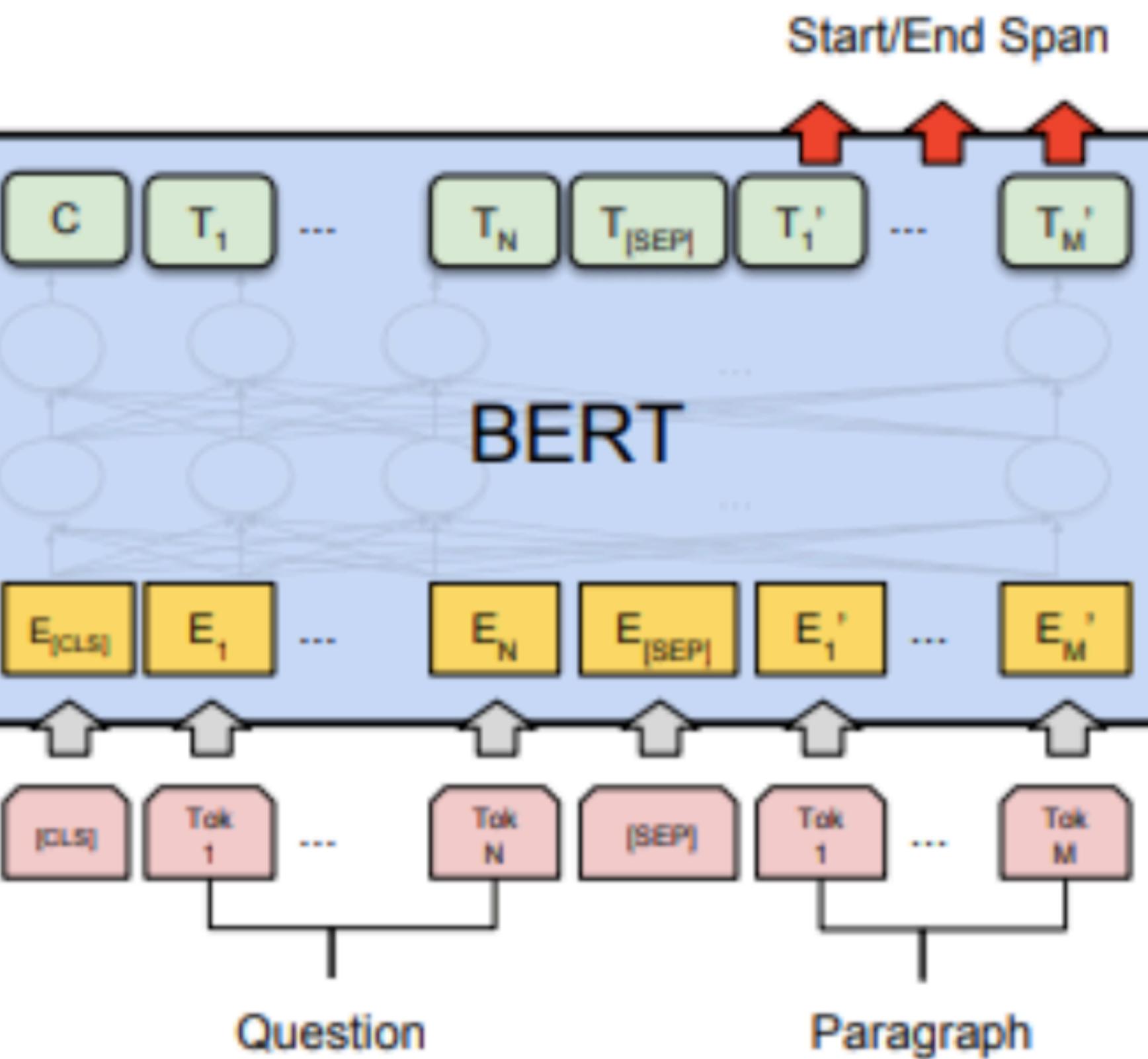
(c) Question Answering Tasks:  
SQuAD v1.1

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

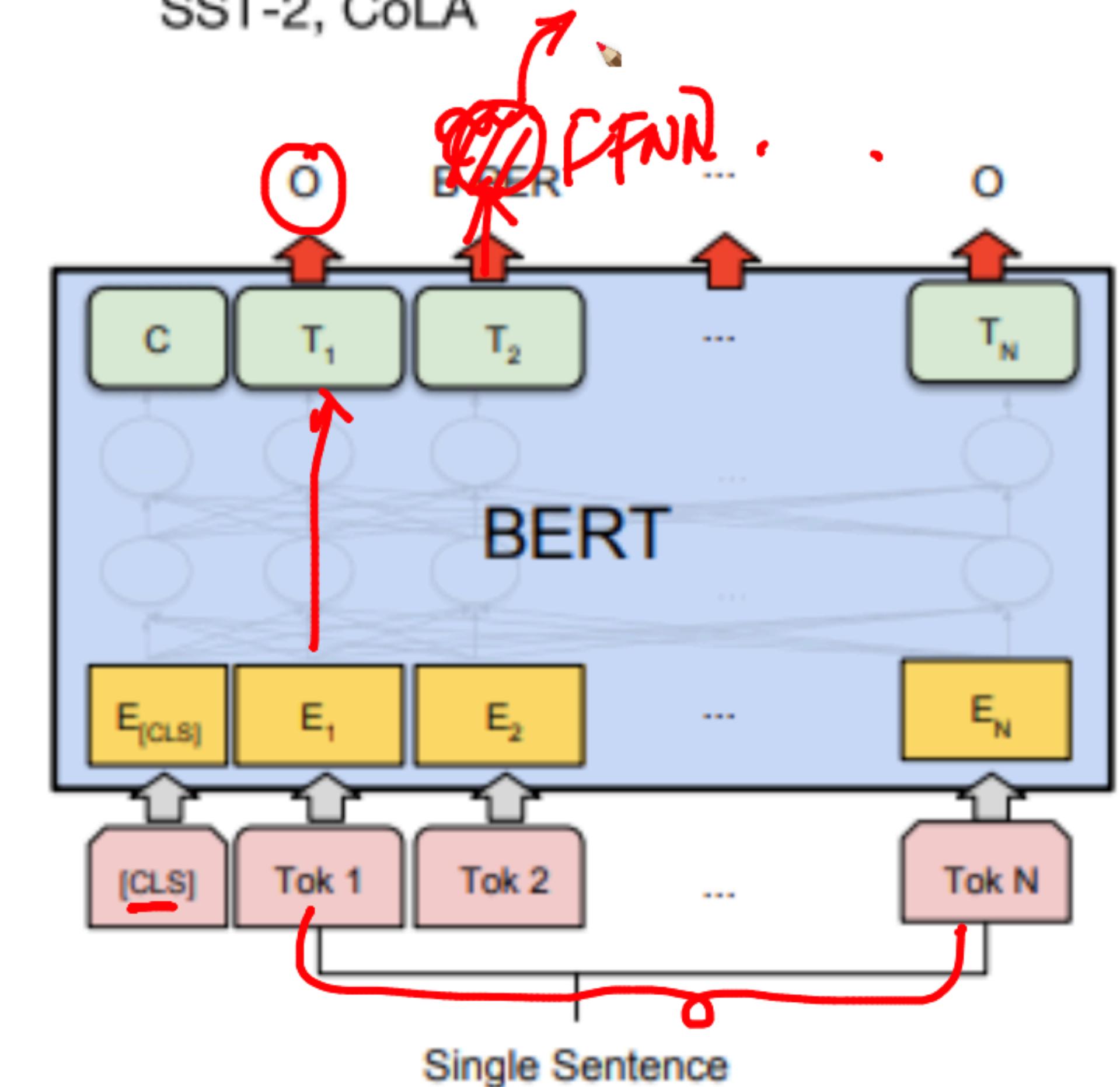


(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

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

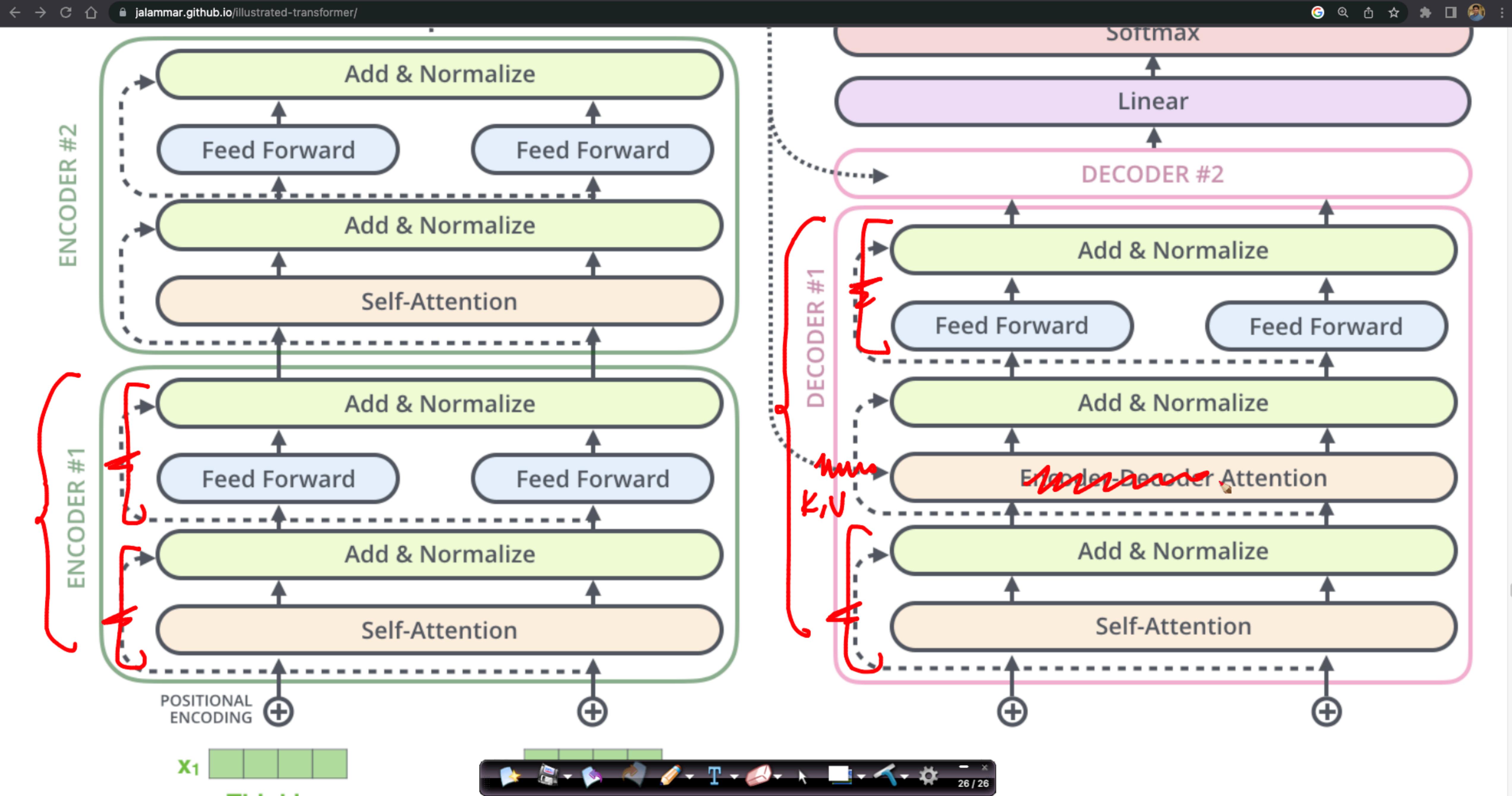


(b) Single Sentence Classification Tasks:  
SST-2, CoLA

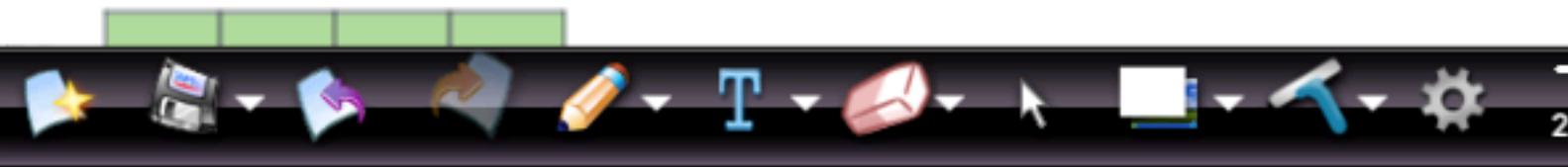


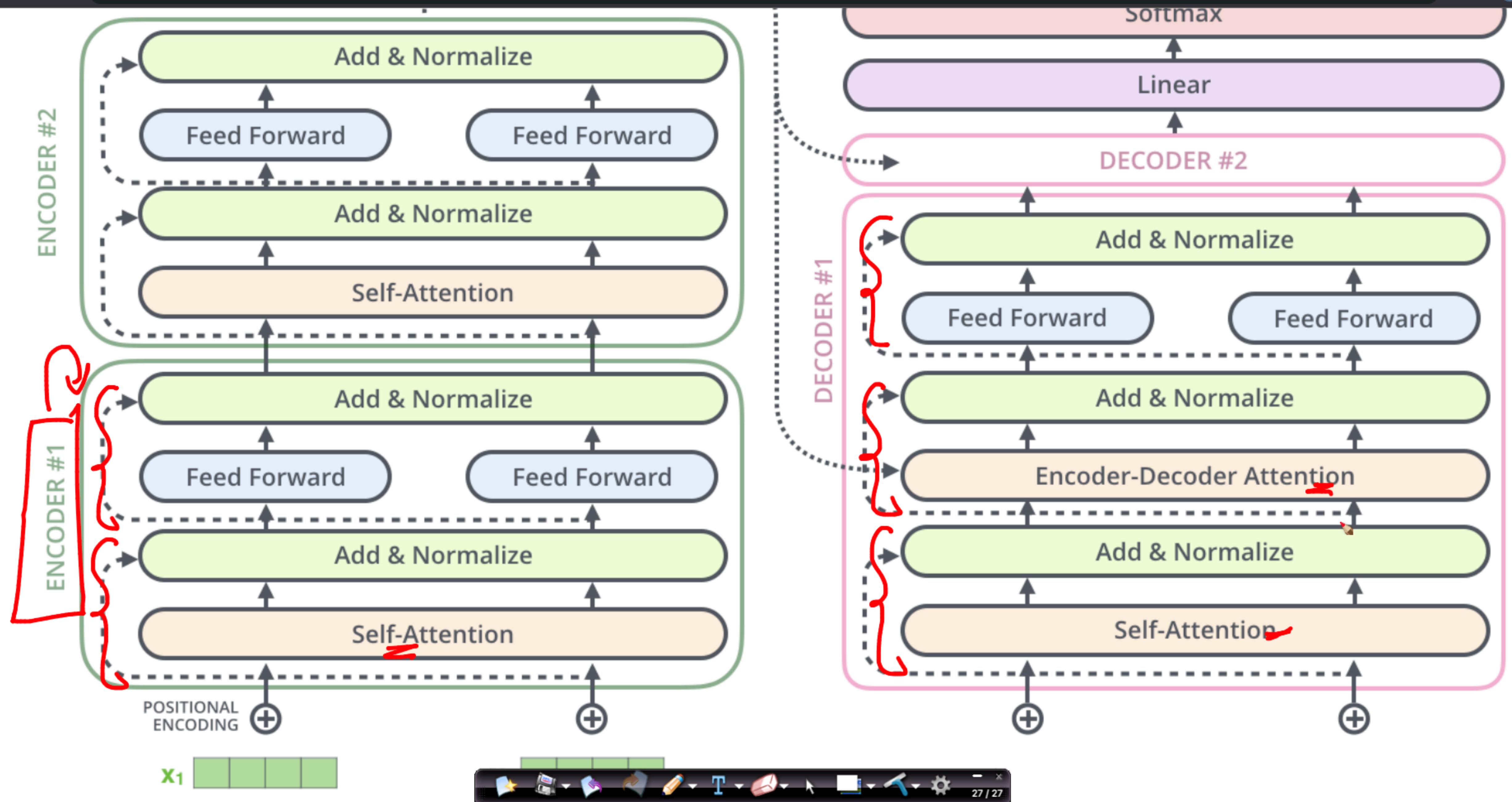
(c) Question Answering Tasks:  
SQuAD v1.1

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

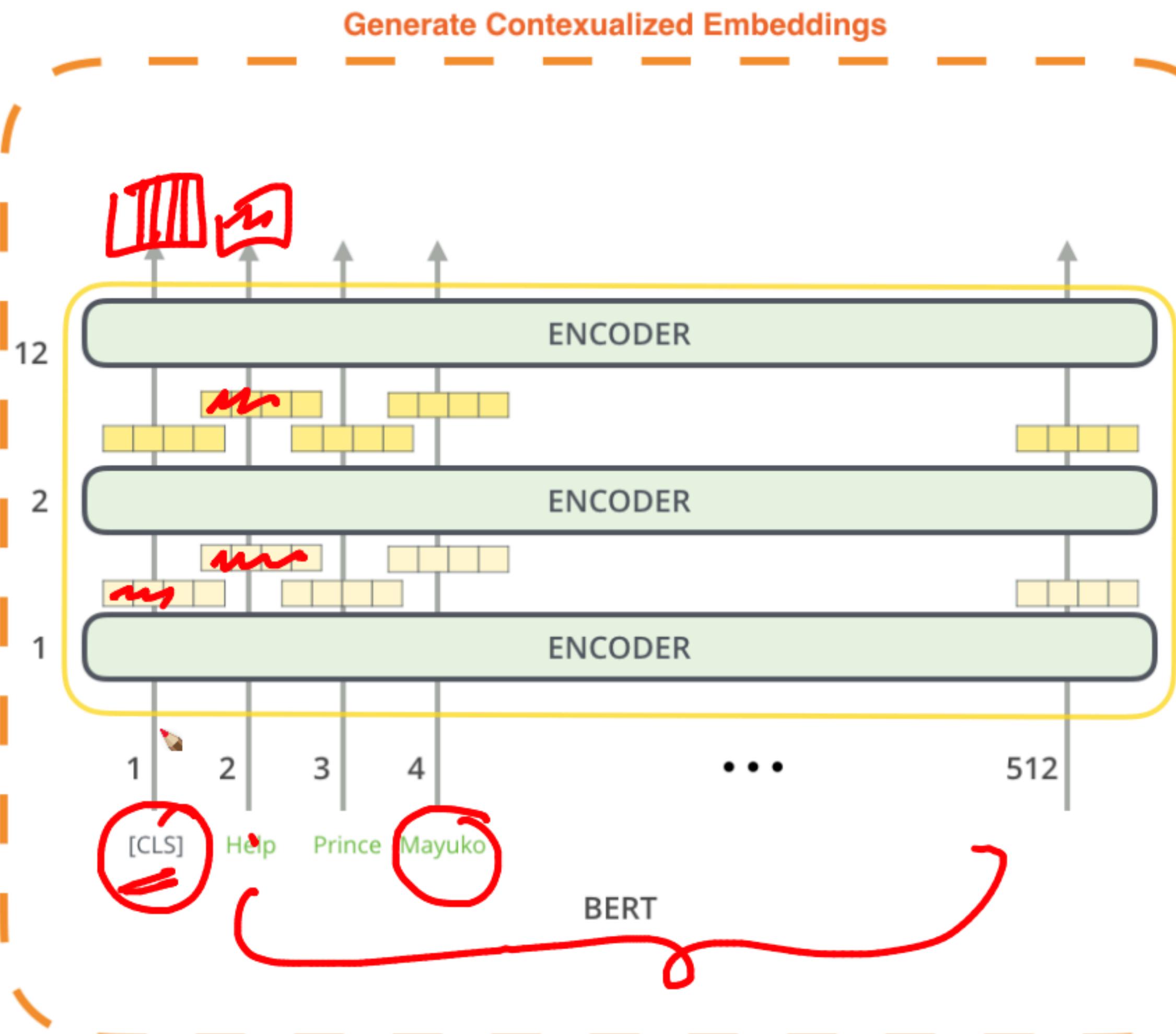


X<sub>1</sub> [ ]

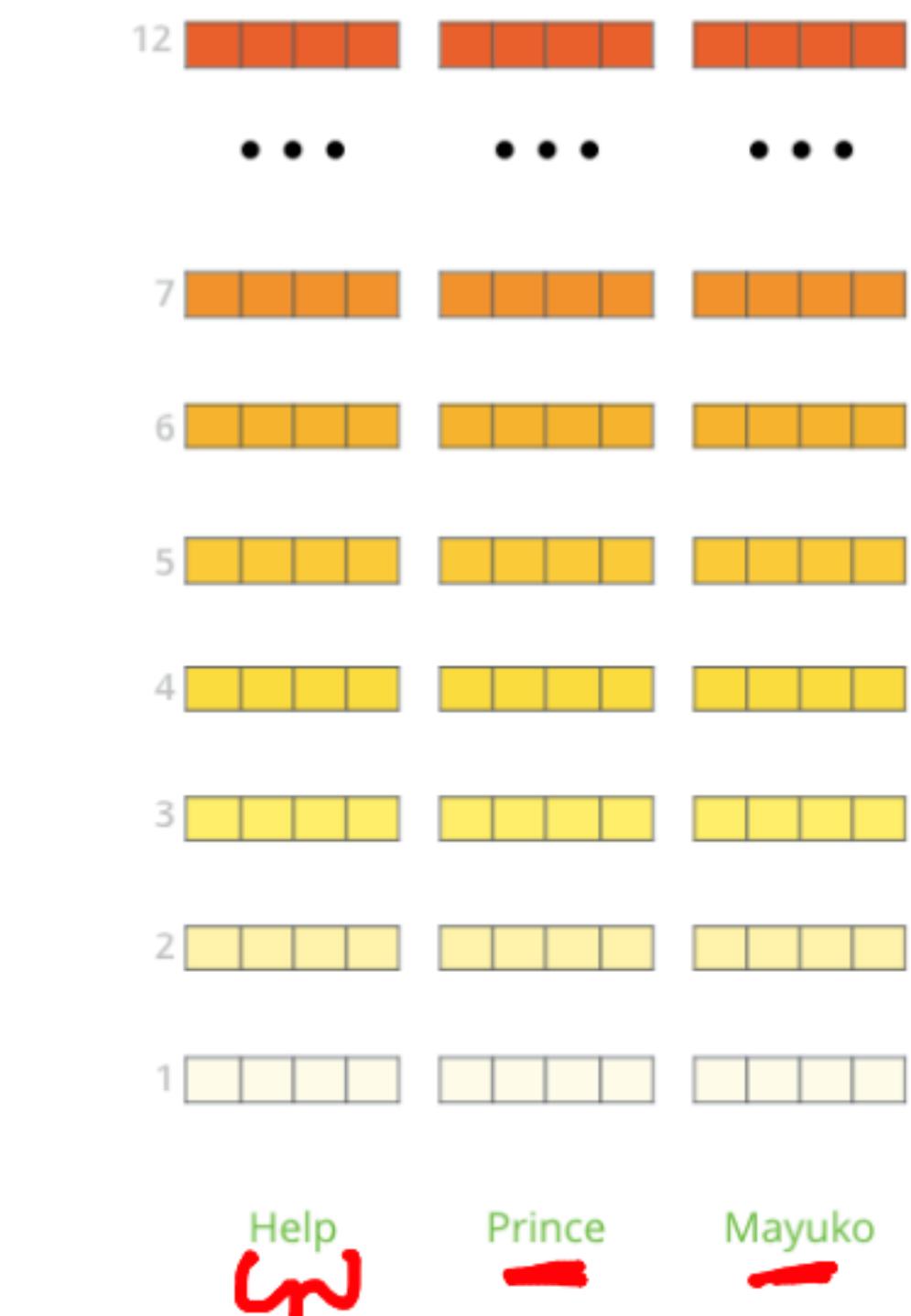




shows yield results not far behind fine-tuning BERT on a task such as named-entity recognition.



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



But which one should we use?

Which vector works best as a contextualized embedding? I would think it depends on the task. The paper examines six choices (Compared to the

What is the best contextualized embedding for “**Help**” in that context  
For named-entity recognition task CoNLL-2003 NER

The diagram illustrates seven different configurations of an NLP architecture, each consisting of an input layer (Embedding) and one or more hidden layers. The configurations are:

- First Layer**: An input embedding layer followed by a single hidden layer (12 units).
- Last Hidden Layer**: An input embedding layer followed by a single hidden layer (12 units).
- Sum All 12 Layers**: An input embedding layer followed by 12 hidden layers (12 units each). A red bracket groups these layers, and a red arrow points from the output to a Dev F1 Score of 95.5.
- Second-to-Last Hidden Layer**: An input embedding layer followed by 11 hidden layers (11 units each).
- Sum Last Four Hidden**: An input embedding layer followed by four hidden layers (12, 11, 10, 9 units respectively). A red bracket groups these layers, and a red arrow points from the output to a Dev F1 Score of 95.9.
- Concat Last Four Hidden**: An input embedding layer followed by four hidden layers (12, 11, 10, 9 units respectively). The outputs of these four layers are concatenated into a single vector (44 units long).

Red annotations include:

- A large red bracket on the left side of the diagram groups the input embedding layer and all hidden layers.
- A red curly brace above the "First Layer" section highlights the input embedding layer.
- A red curly brace above the "Sum All 12 Layers" section highlights the 12 hidden layers.
- A red curly brace above the "Sum Last Four Hidden" section highlights the four hidden layers.
- A red arrow points from the "Sum All 12 Layers" output to a circled Dev F1 Score of 95.5.
- A red arrow points from the "Sum Last Four Hidden" output to a circled Dev F1 Score of 95.9.
- A red arrow points from the "Concat Last Four Hidden" output to a circled Dev F1 Score of 96.1.
- A red checkmark is placed next to the "First Layer" configuration.
- A red checkmark is placed next to the "Concat Last Four Hidden" configuration.
- A red circle highlights the Dev F1 Score 91.0.
- A red circle highlights the Dev F1 Score 94.9.
- A red circle highlights the Dev F1 Score 96.1.
- A red circle highlights the Dev F1 Score 95.5.
- A red circle highlights the Dev F1 Score 95.6.

(Q)

{ Train  $\Rightarrow$   
10K sentencesclassfn: 3-class  
clsBERT  
(512)~ few million params

(?)

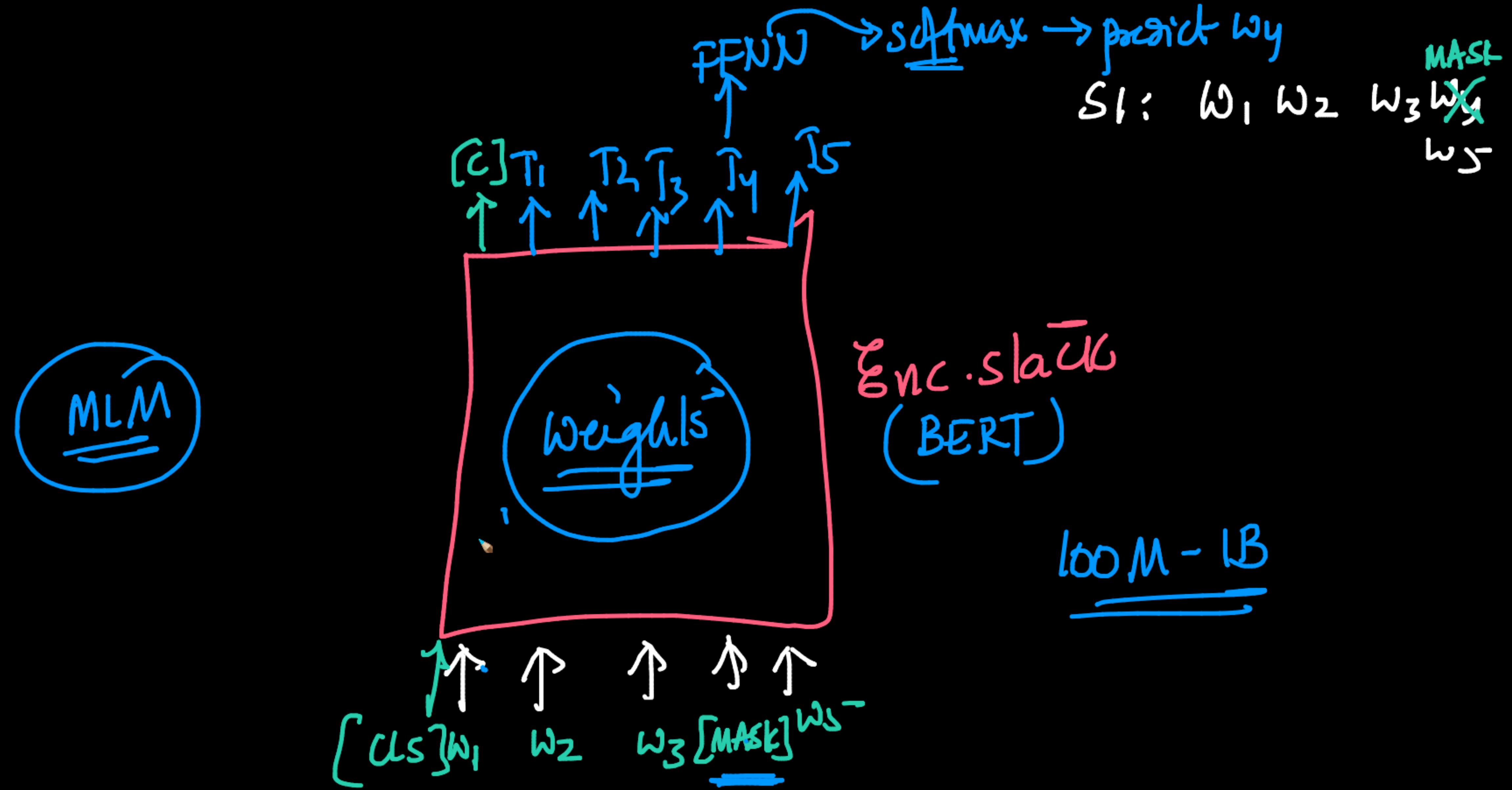
fine tune a pre-trained  
modelCNNs → ?  
↓

ImageNet

label  
10-14 MB

English- sentences →  ✓  
 $w_1 w_2 w_4 w_b w_f$

BERT → Masked language models  
→ self-supervised learning



(Q)

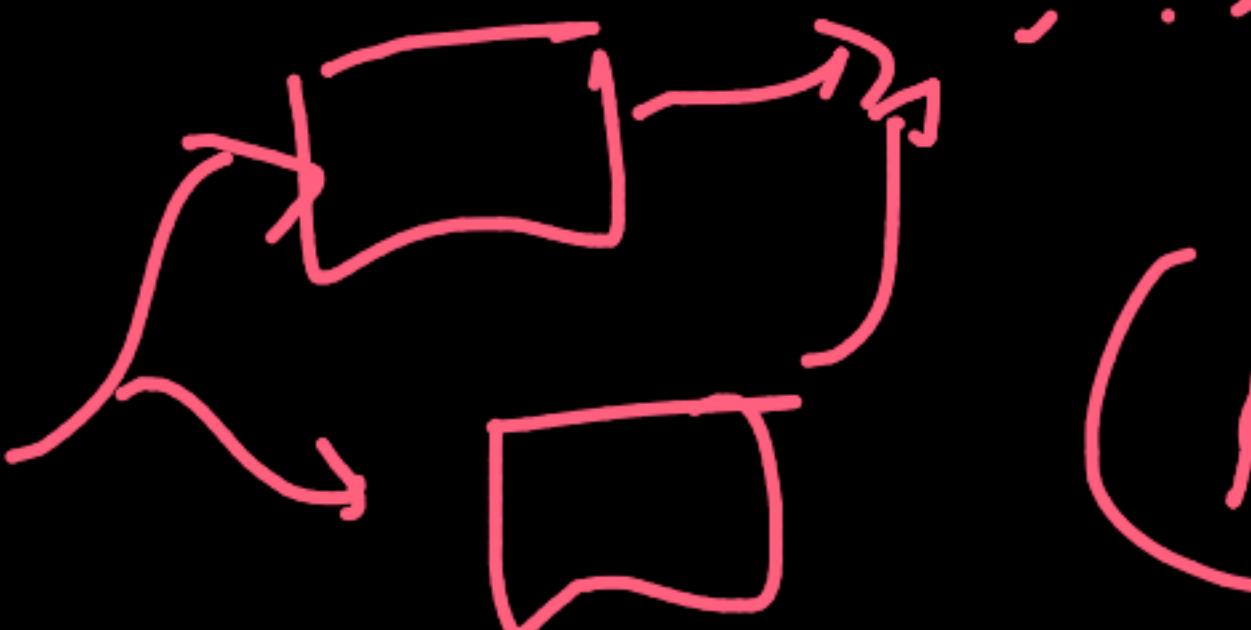
Bert → pre-trained + fine-tune  
(MLM)

@ runtime

low-latency

BERT

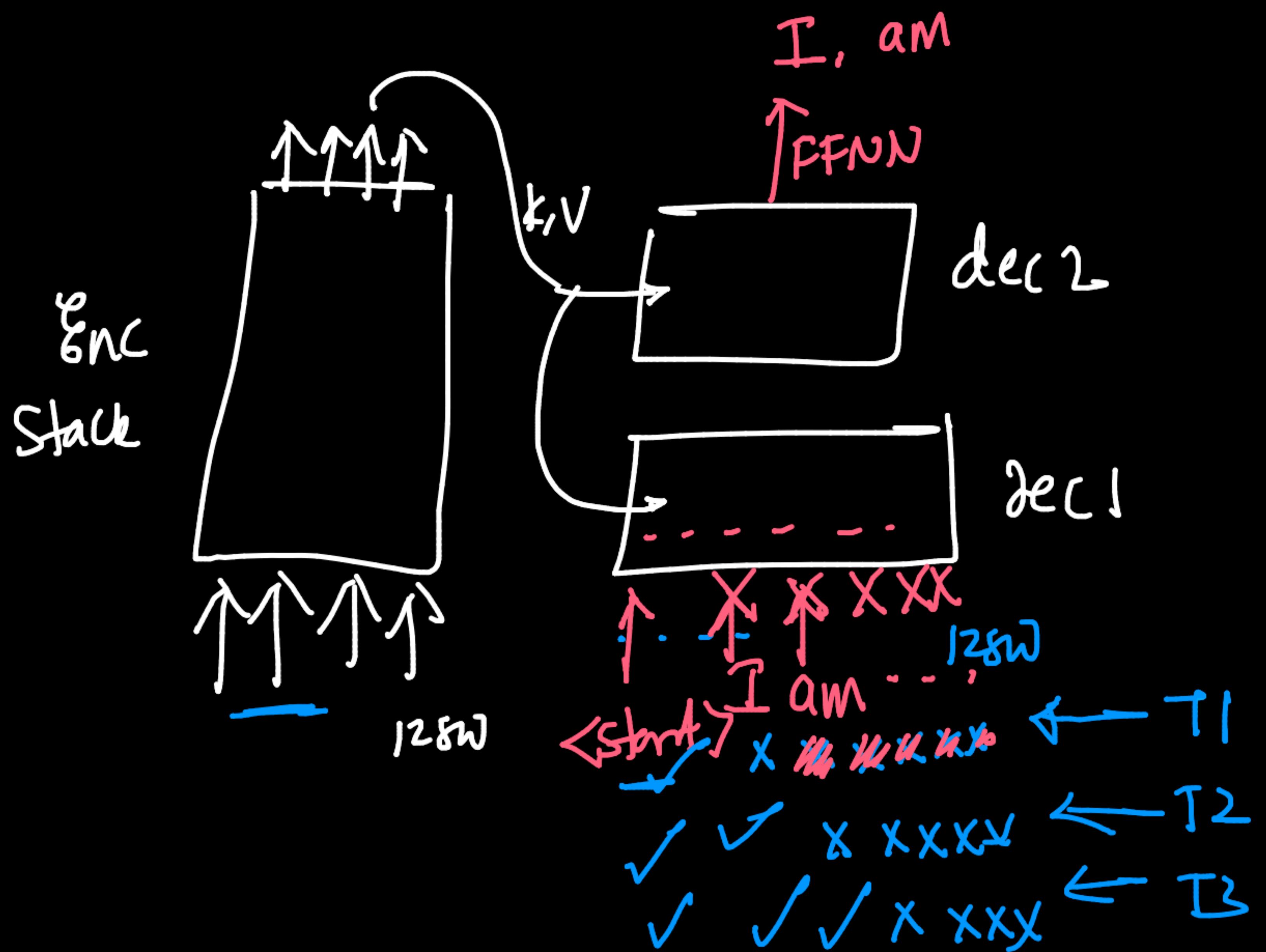
DistilBERT → ✓  
which performs well



(Knowledge Distill)

HINT:

Large  
→ small



+ Code + Text

- learning rate
  - number of epochs

```
from transformers import create_optimizer  
import tensorflow as tf
```

```
# The number of training steps is the number of samples in the dataset, divided by the batch size then multiplied  
# by the total number of epochs. Note that the tf_train_dataset here is a batched tf.data.Dataset,  
# not the original Hugging Face Dataset, so its len() is already num_samples // batch_size.
```

num epochs = 3

```
num_train_steps = len(tf_train_dataset) * num_epochs
```

```
optimizer, schedule = create_optimizer(
```

init\_lr=5e-5,

`num_warmup_steps=0,`

```
num_train_steps=num_train_steps,
```

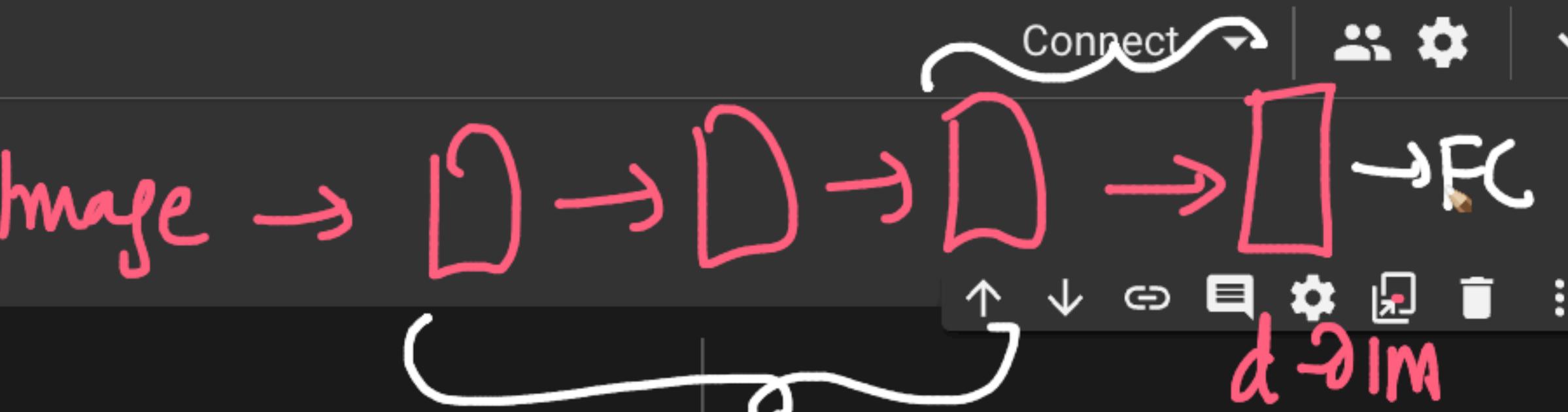
`weight_decay_rate=0.01,`

)

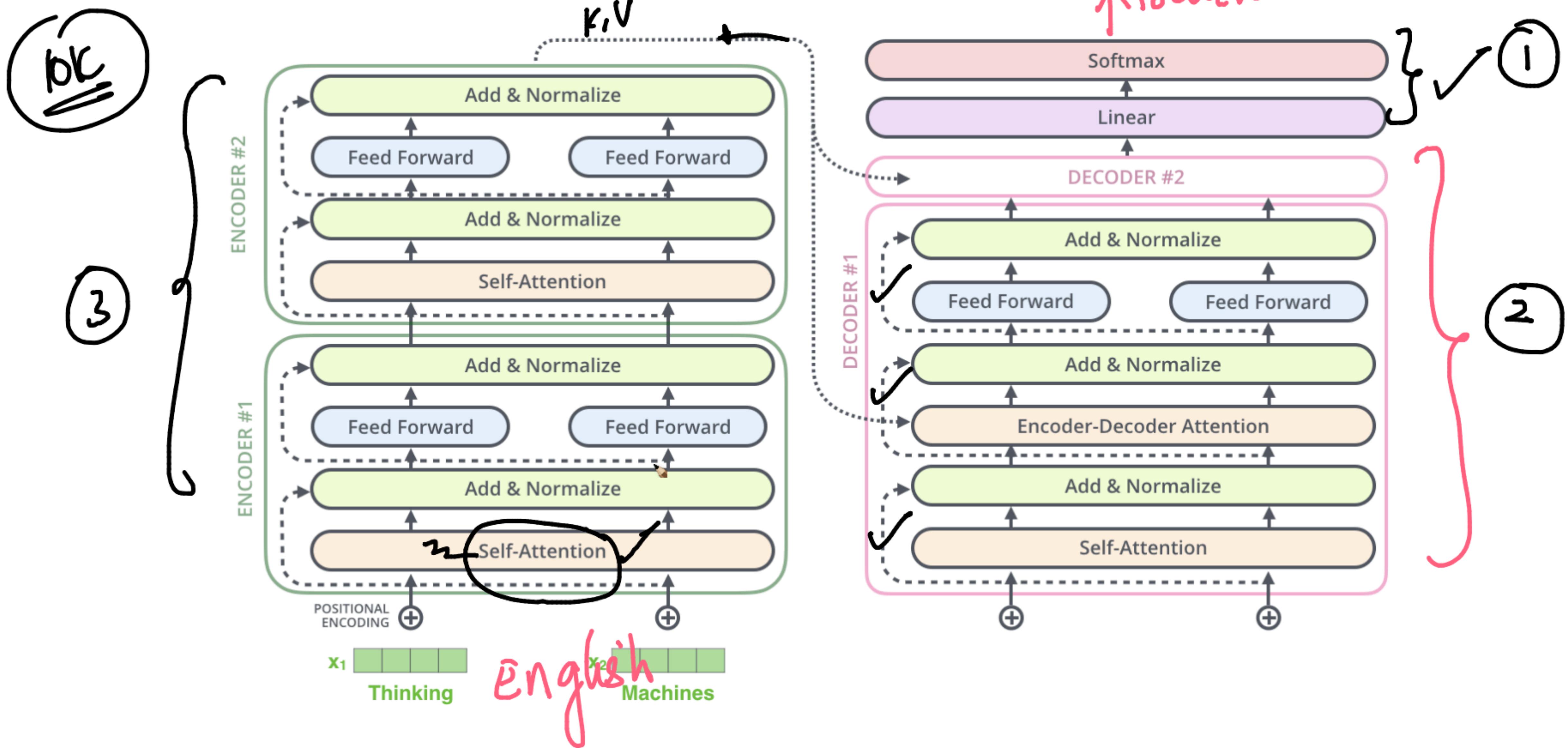
```
model.compile(optimizer=optimizer)
```

No loss specified in `compile()` - the model's internal loss computation will be used as the loss. Don't panic - this

[ ]



decoders, it would look something like this:



+ Code + Text

Connect ▾



- labels



```
from transformers import DataCollatorForSeq2Seq
from transformers import TFAutoModelForSeq2SeqLM, Seq2SeqTrainingArguments, Seq2SeqTrainer

model = TFAutoModelForSeq2SeqLM.from_pretrained(model_checkpoint, from_pt=True)

data_collator = DataCollatorForSeq2Seq(tokenizer=tokenizer, model=model, return_tensors="tf")
```

All PyTorch model weights were used when initializing TFMarianMTModel.

All the weights of TFMarianMTModel were initialized from the PyTorch model.

If your task is similar to the task the model of the checkpoint was trained on, you can already use `TFMarianMTModel`.

```
[ ] # Lets see content of data collator
```

```
batch = data_collator([tokenized_datasets["train"][i] for i in range(1, 3)])
batch.keys()
```

```
dict_keys(['input_ids', 'attention_mask', 'labels', 'decoder_input_ids'])
```

```
# Lets see how decoder input is mapped to id  
batch["decoder input ids"]
```