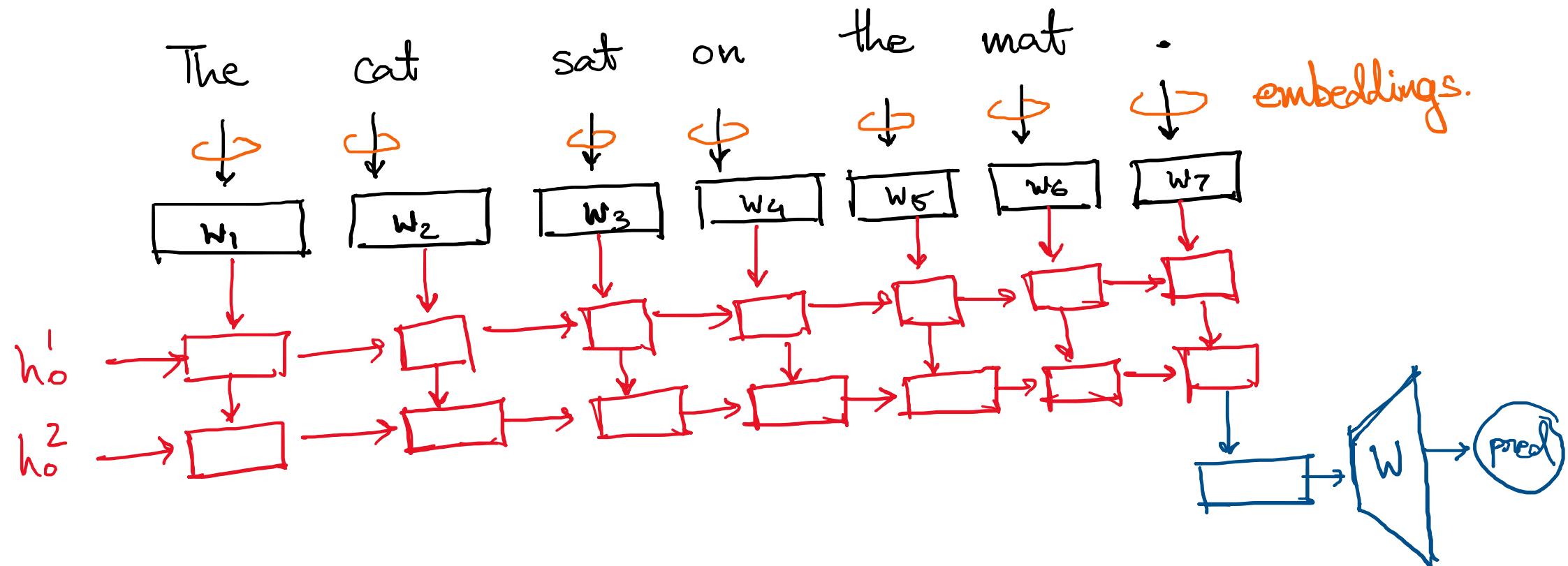


BERT & other Transformer-
based language models: an overview

Soham Pal.

27th May 2020

The need for embeddings



Word embeddings must capture:

① Complex characteristics of word use

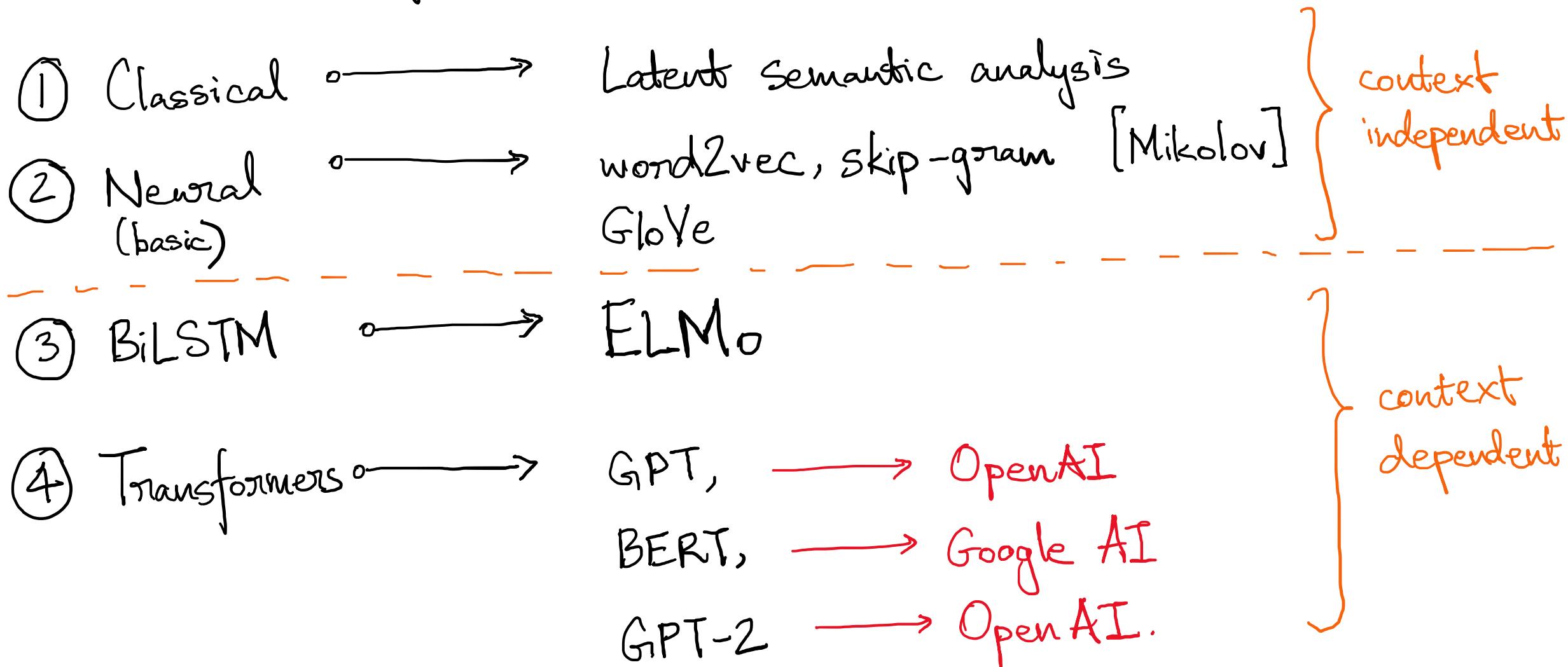
→ syntax

→ Semantics

② How these uses vary across linguistic contexts

→ polysemy.

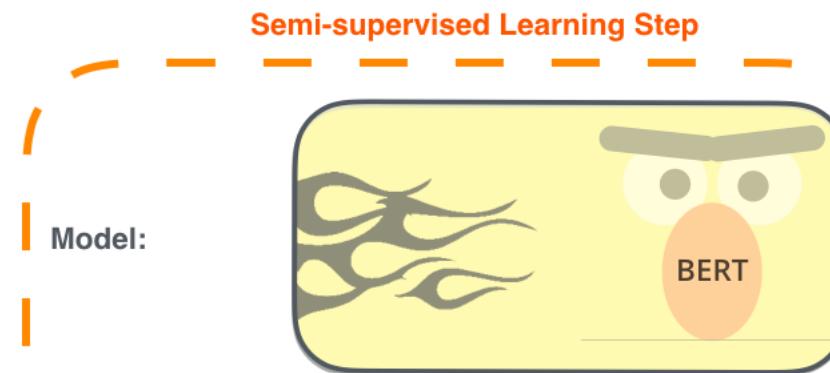
Word Embeddings: a timeline



Overall framework

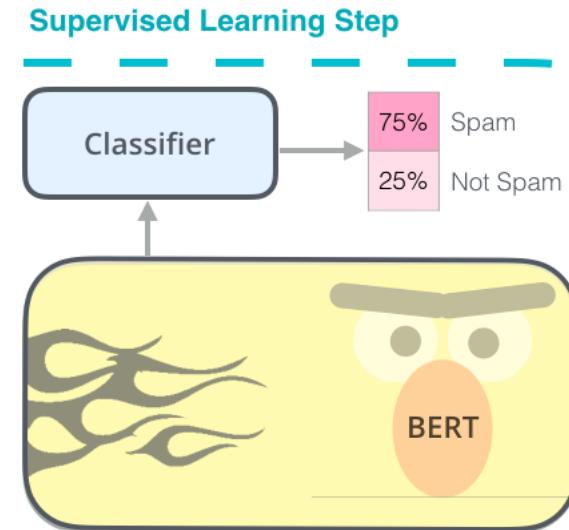
1 - Semi-supervised training on large amounts of text (books, wikipedia..etc).

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.



Model: BERT
Dataset: Books, Wikipedia
Objective: Predict the masked word (language modeling)

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



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

Approaches to the

Supervised
Learning Step.

Fine-tuning-based

BERT, OpenAI GPT-2, etc.

all network parameters
retrained.

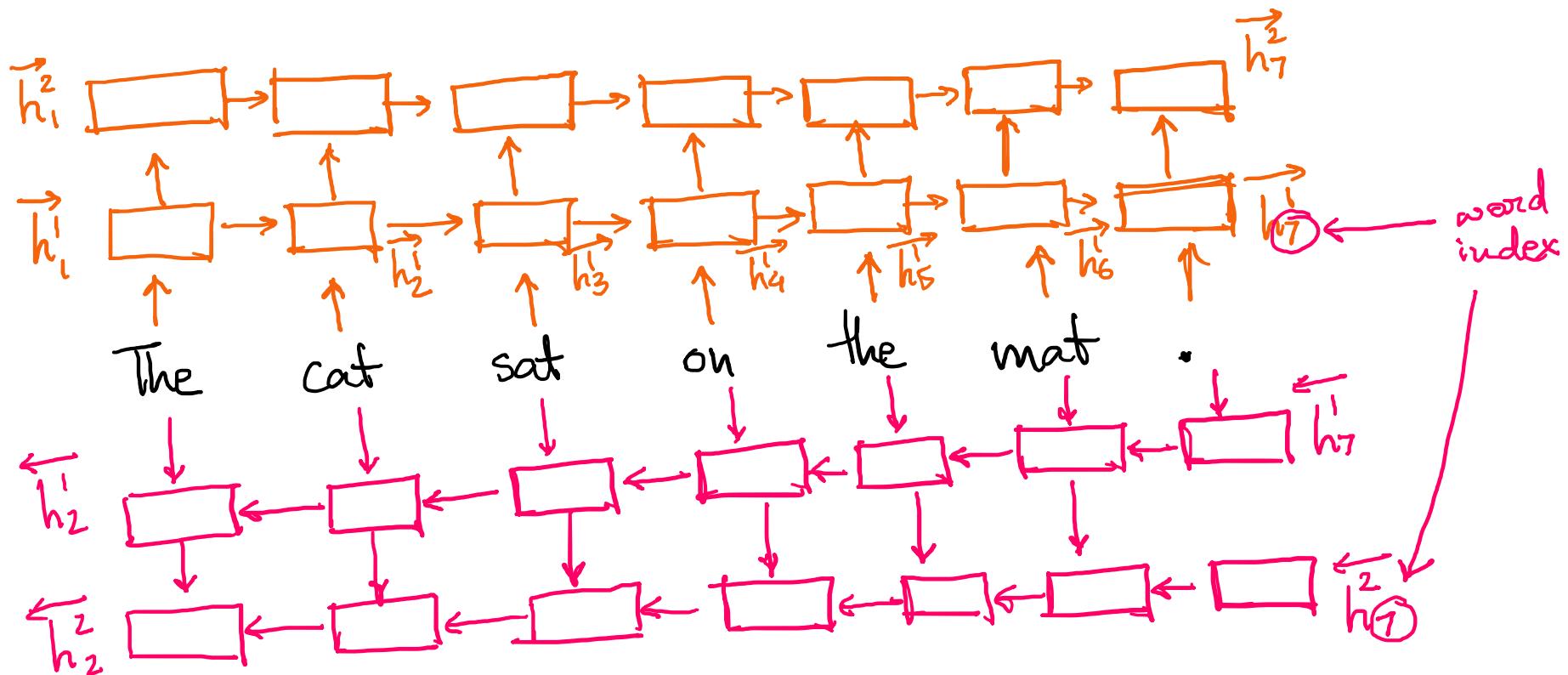
Feature-based.

ELMo.

only task-specific
parameters
learned

ELMo

- start with large dataset (corpus: monolingual, unannotated)
- train biLSTM language model.



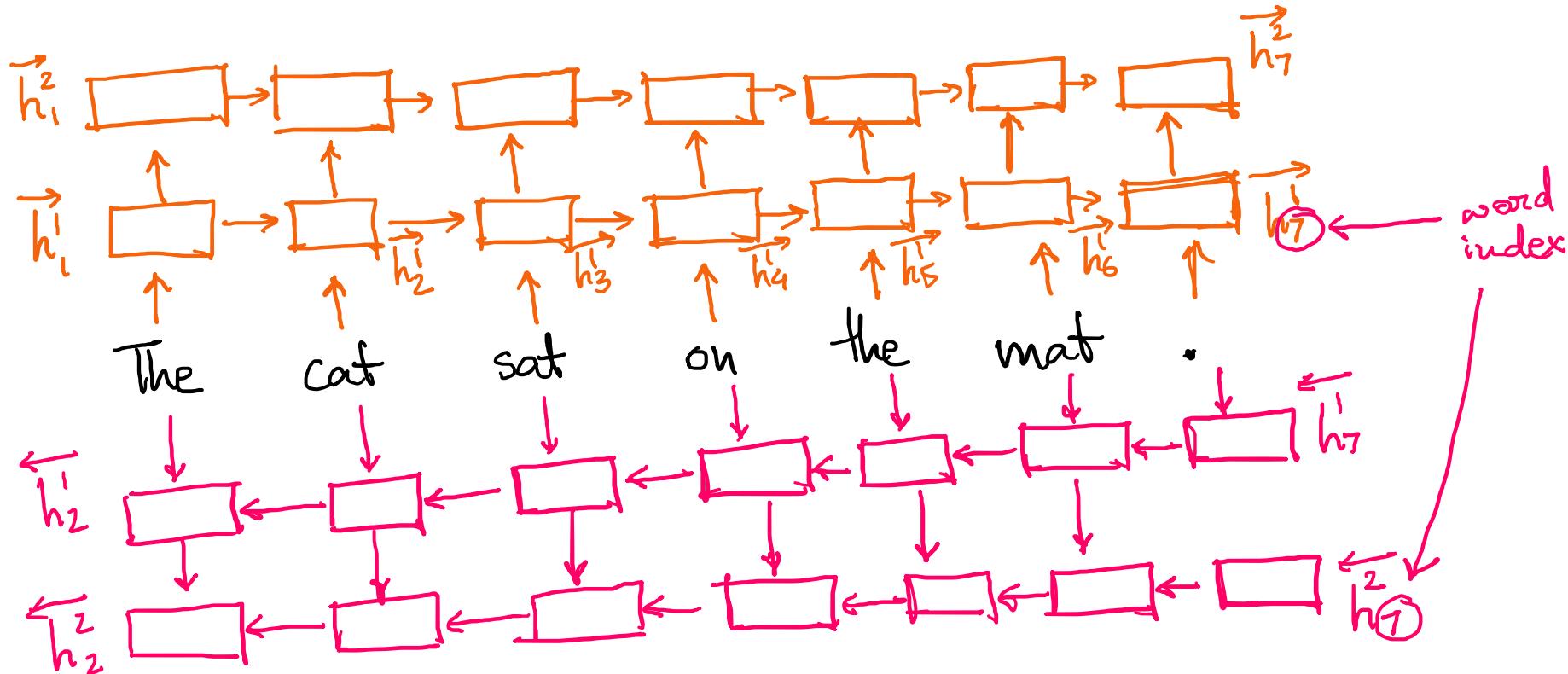
ELMo (cont'd).

$$\text{ELMO}_k = \lambda \sum_{j=1}^2 s_j [\overrightarrow{h_k^j}, \overleftarrow{h_k^j}]$$

word index.

sum across layers.

forward backward.



ELMO (contd).

$$\text{ELMO}_k = \lambda \left[\sum_{j=1}^2 s_j \begin{bmatrix} \vec{h}_k^j & \vec{h}_k^j \end{bmatrix} \right]$$

(End of pretraining)

forward propagation ONLY

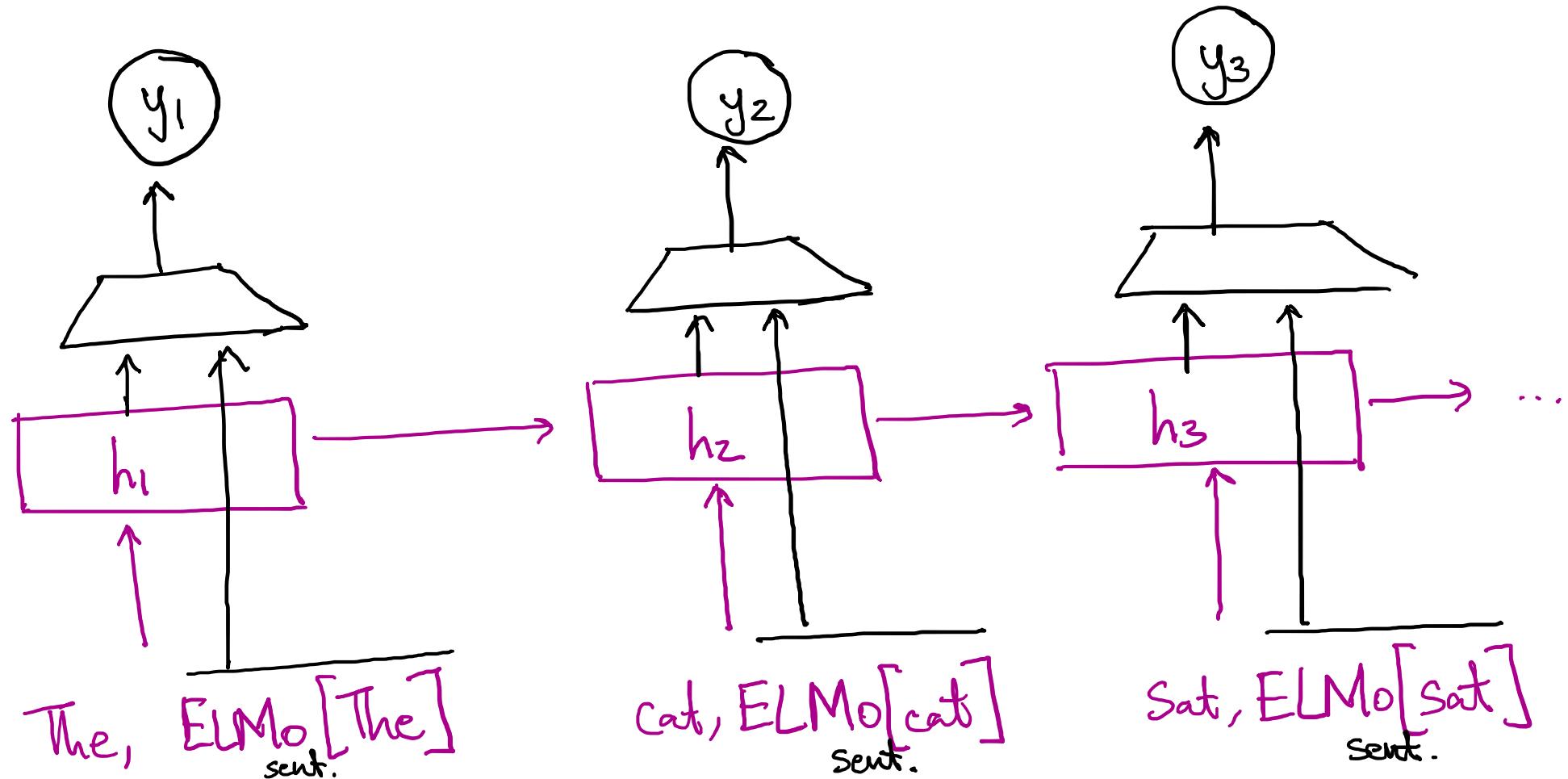
backward propagation ONLY

\otimes

The diagram illustrates the computation of the ELMO vector. It starts with the formula $\text{ELMO}_k = \lambda \left[\sum_{j=1}^2 s_j \begin{bmatrix} \vec{h}_k^j & \vec{h}_k^j \end{bmatrix} \right]$. Below the formula, a timeline is shown with a horizontal axis. A bracket labeled "forward propagation ONLY" spans from left to right, and another bracket labeled "backward propagation ONLY" spans from right to left, both ending at a circled "X". This visualizes that information only flows forward or backward through the network, never both ways simultaneously.

No cross
propagation
of information.

ELMO: Task-specific training -



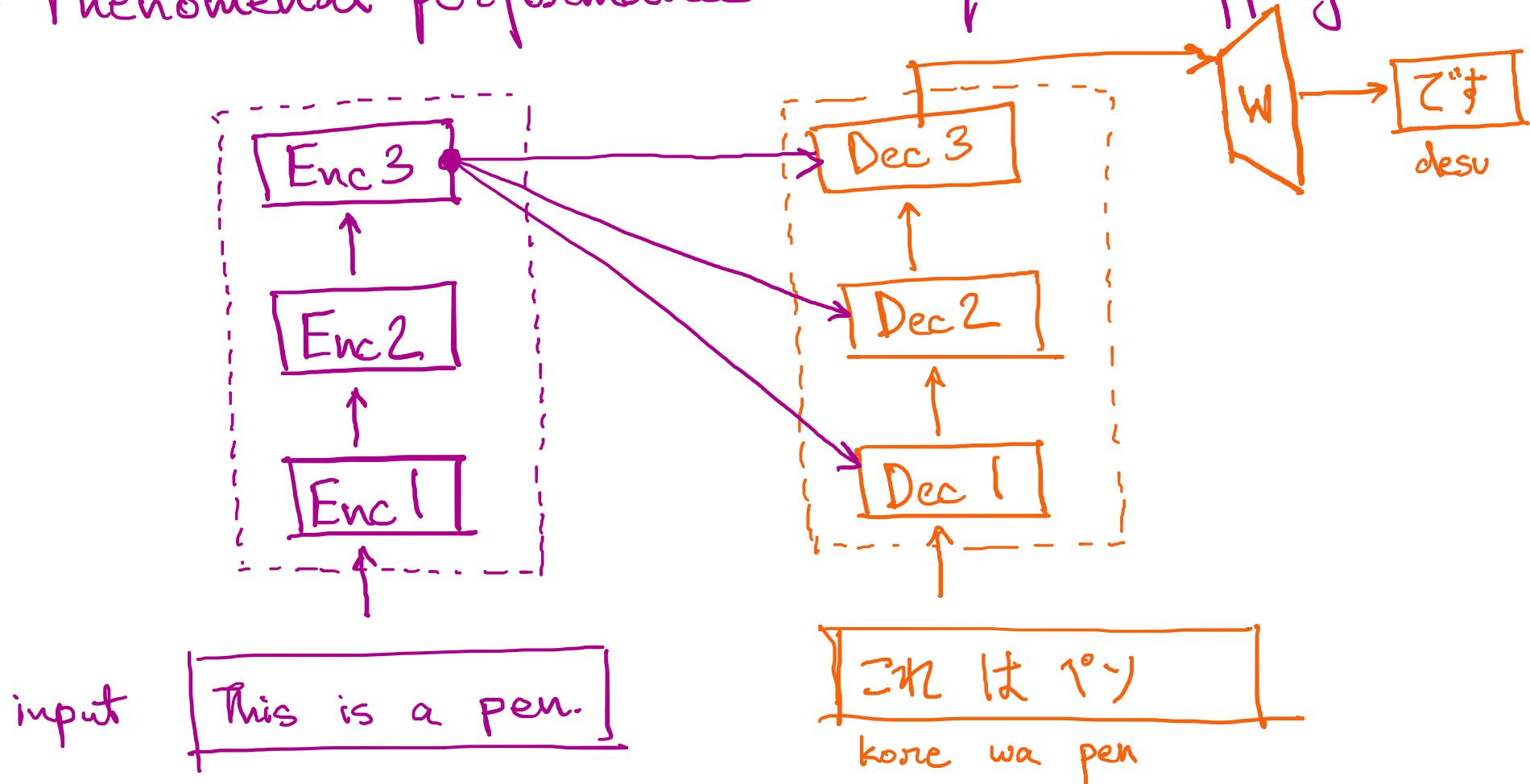
Freeze LM parameters, learn only λ and s .

Problems with ELMo:

- basic architecture, trained on limited data.
"Fixed" by GPT.
- no mixing of information from both directions.
"Fixed" by BERT.

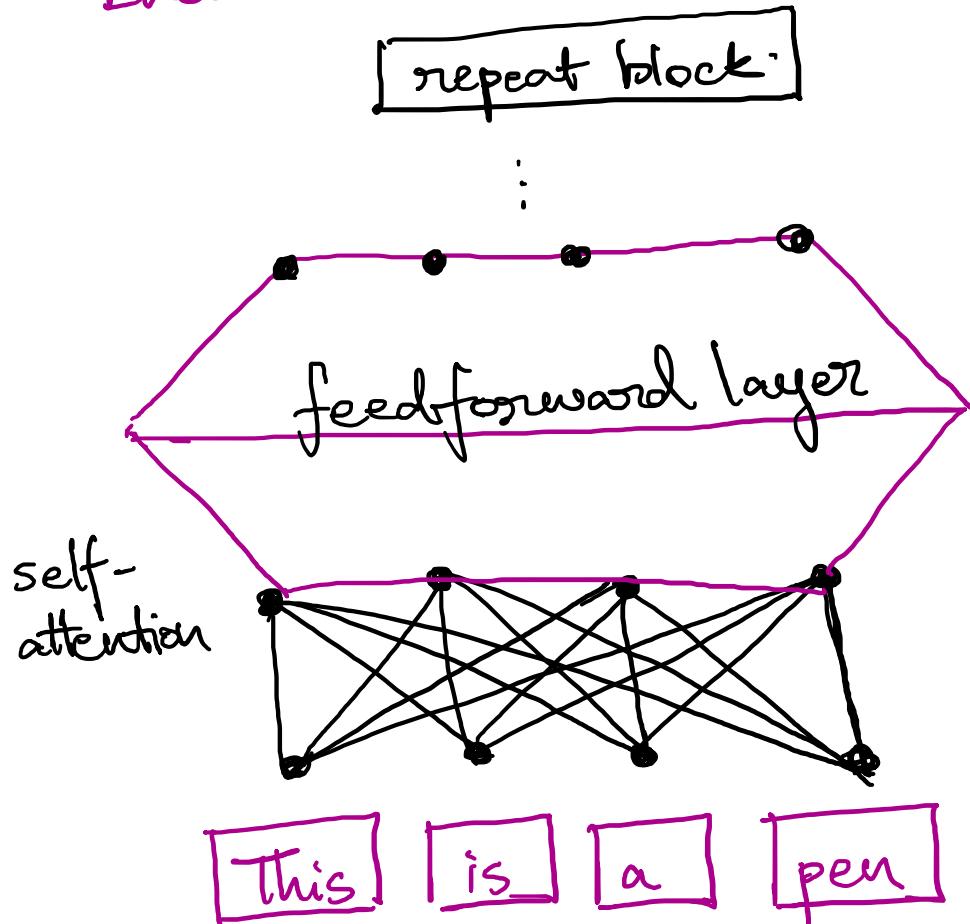
Recap: Transformers

→ Phenomenal performance on sequence-mapping tasks.



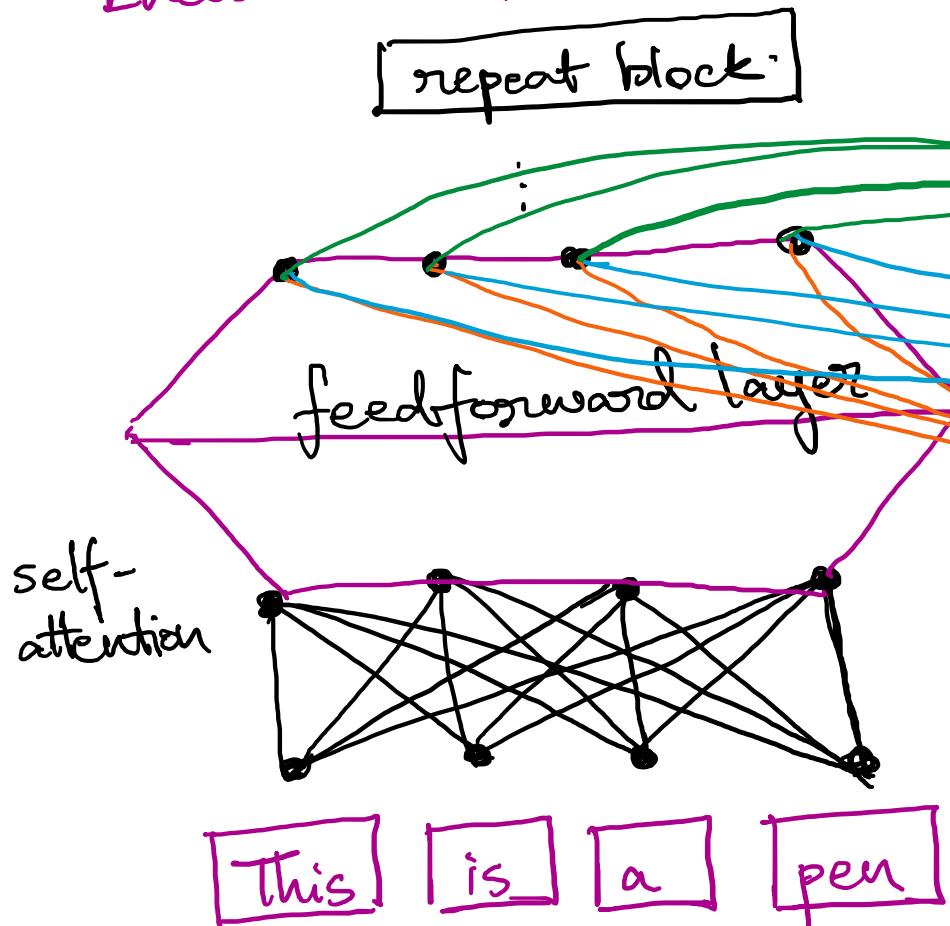
Recap: Transformers (cont'd)

Encoder block

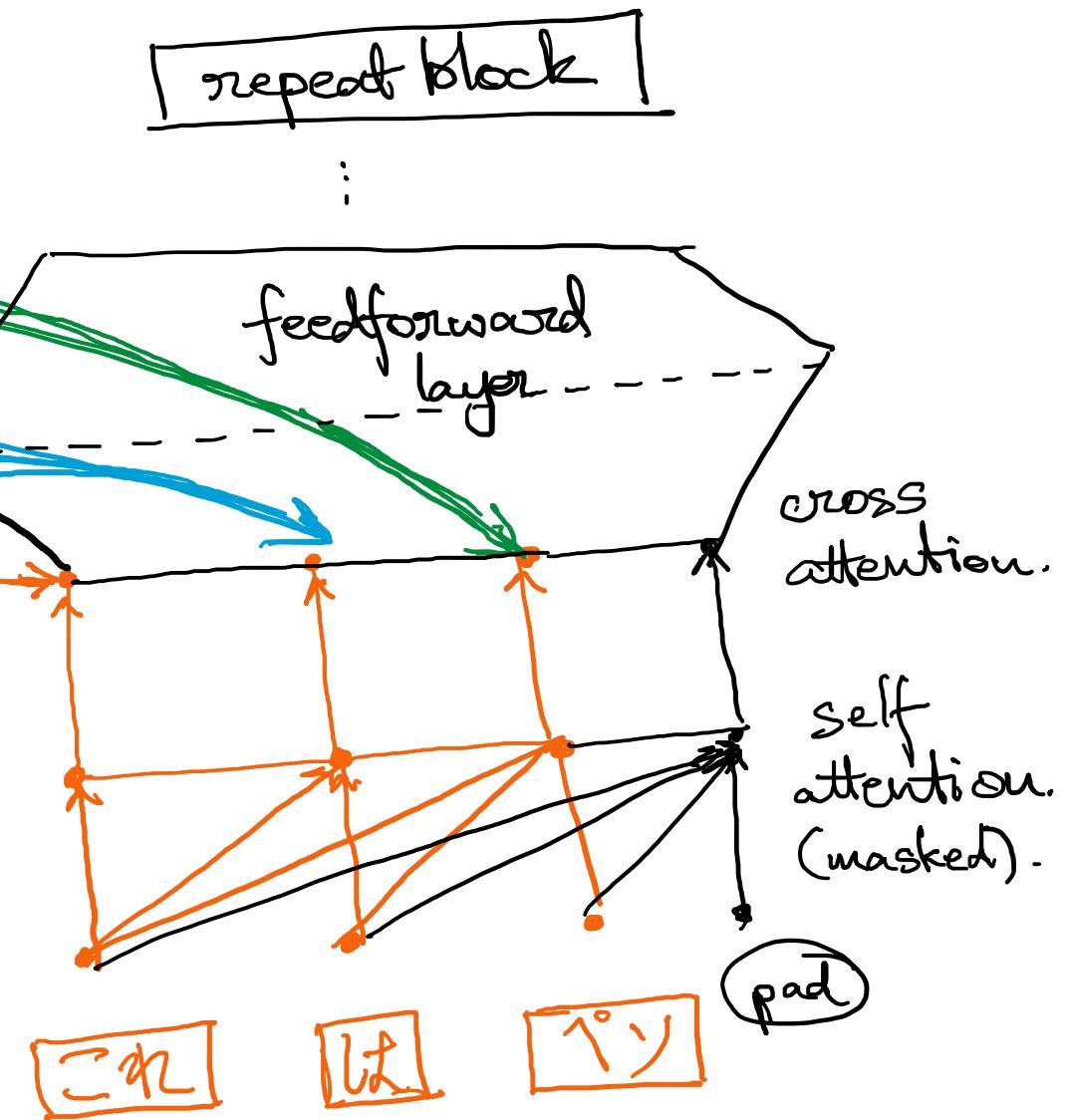


Recap: Transformers (cont'd)

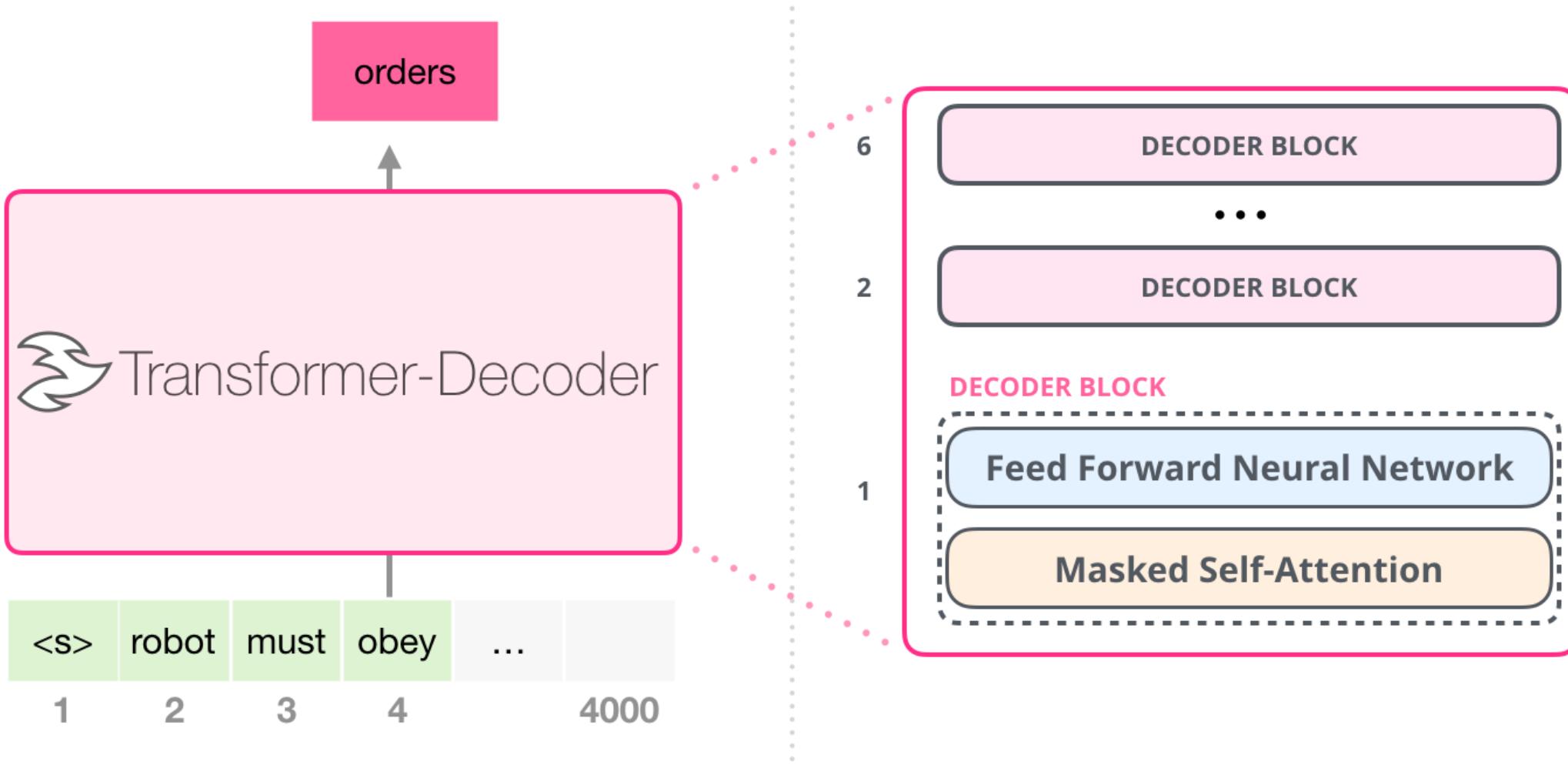
Encoder block



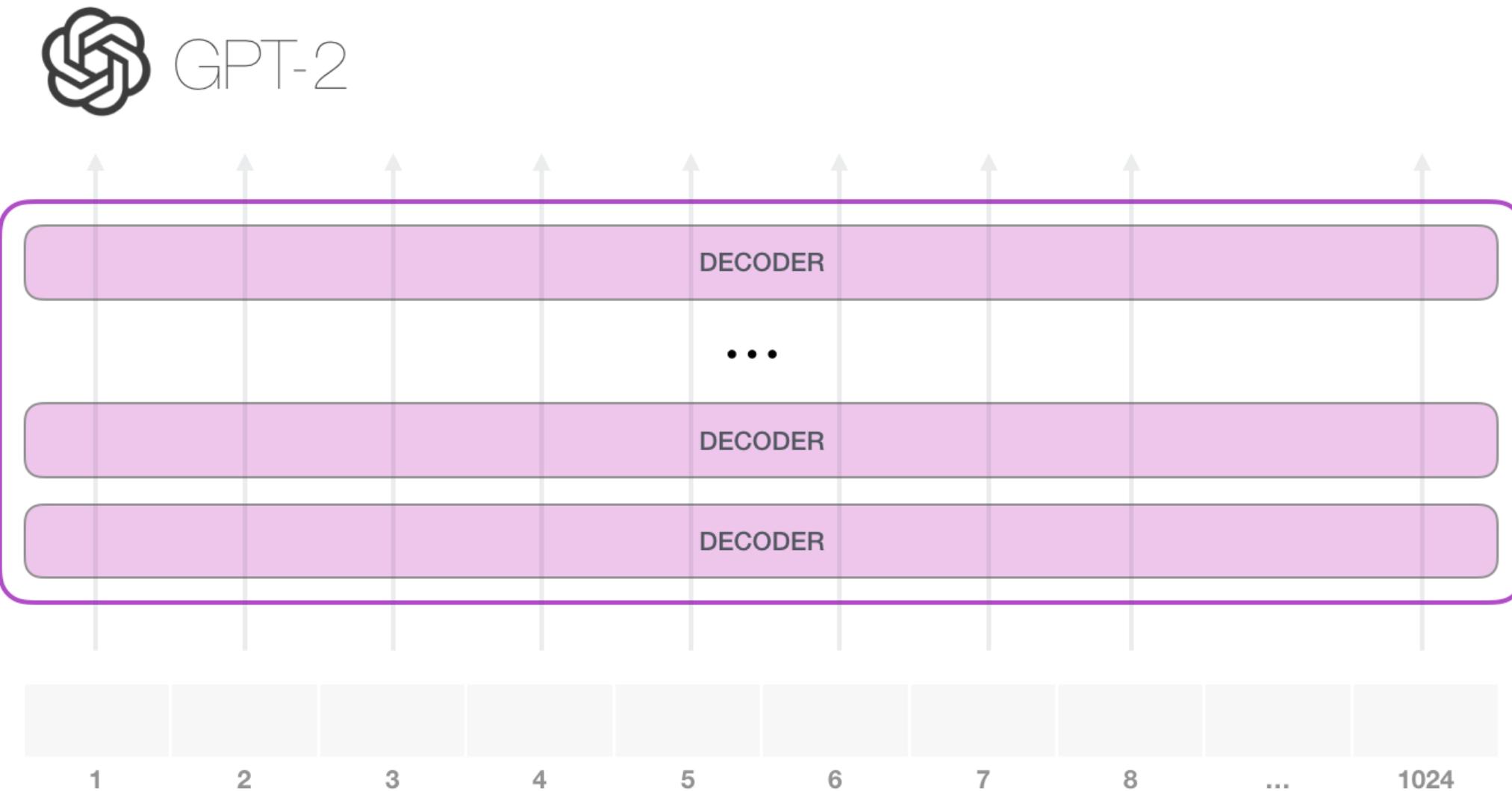
Decoder block



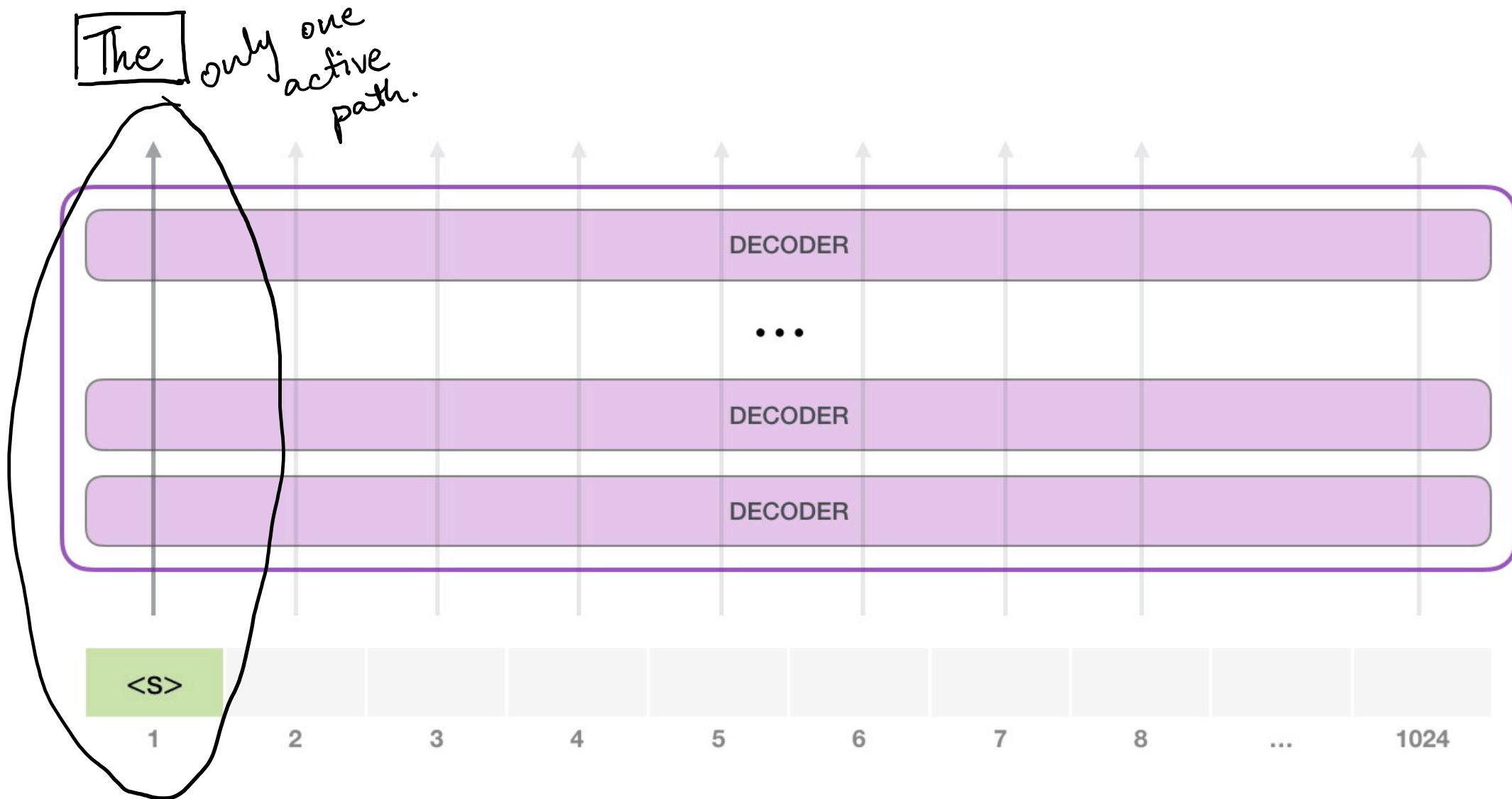
The new: TRANSFORMER-DECODER BLOCK.



The new: TRANSFORMER-DECODER BLOCK.

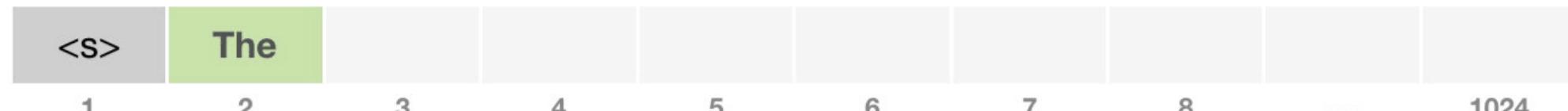
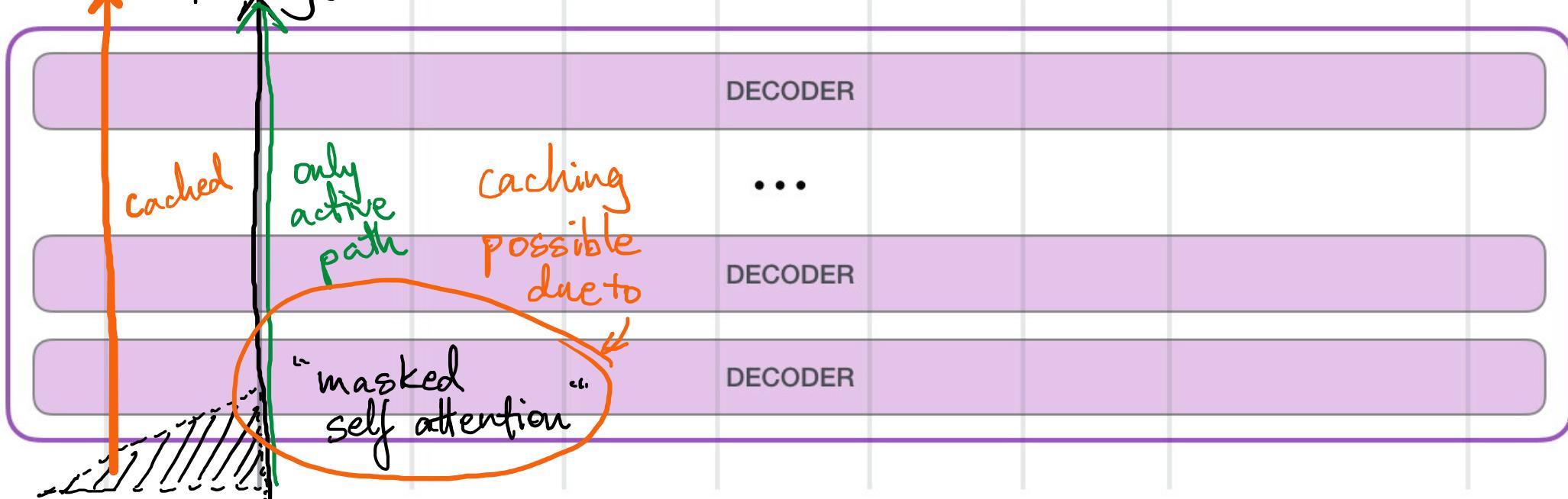


The new: TRANSFORMER-DECODER BLOCK.

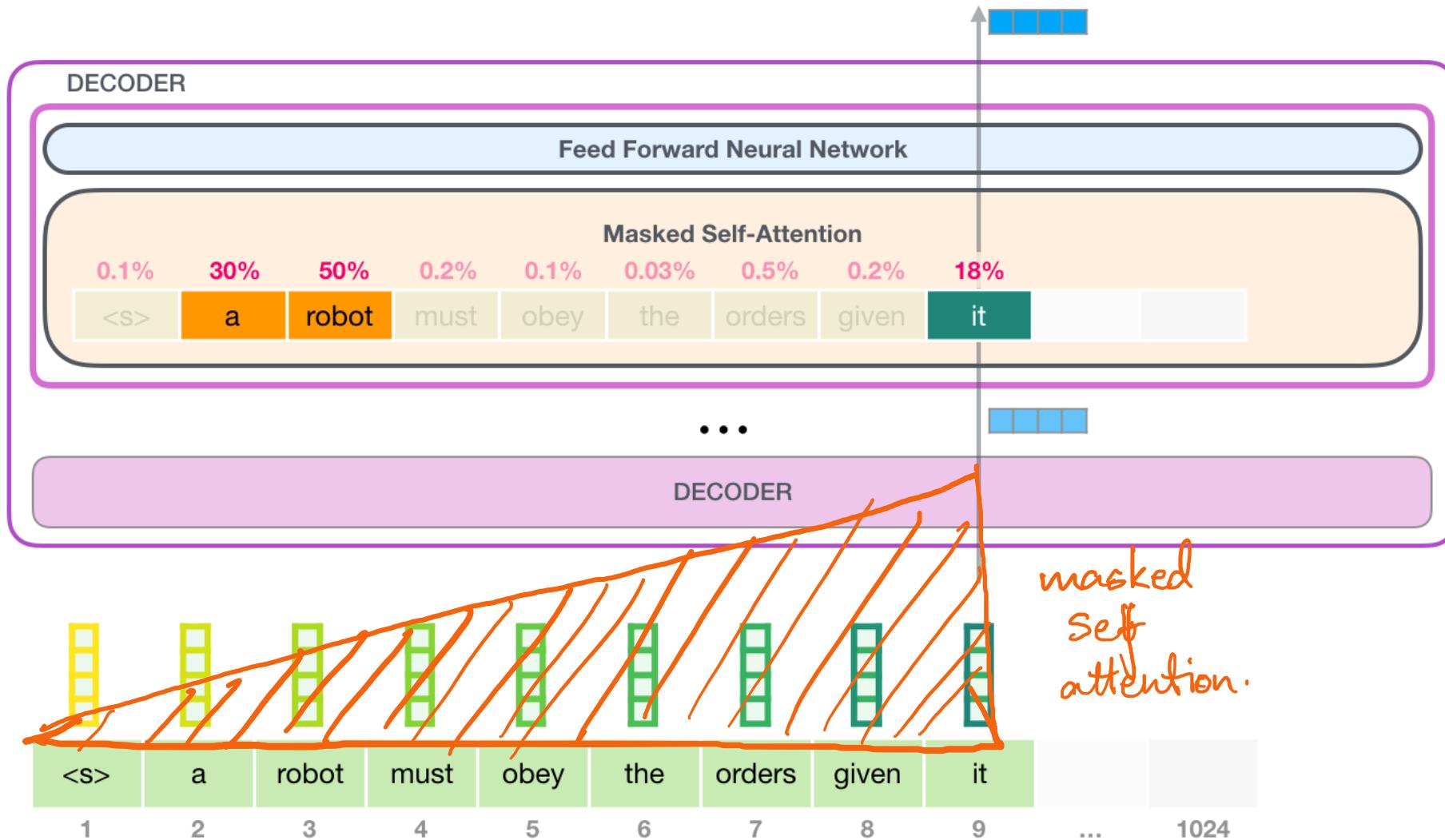


The new: TRANSFORMER-DECODER BLOCK.

Each layer of GPT-2 has retained its own "interpretation" of the first token & will use it in processing the second token.



Peeling back the attention mechanism:



GPT: Putting it all together.

Transformer-decoder we just saw.

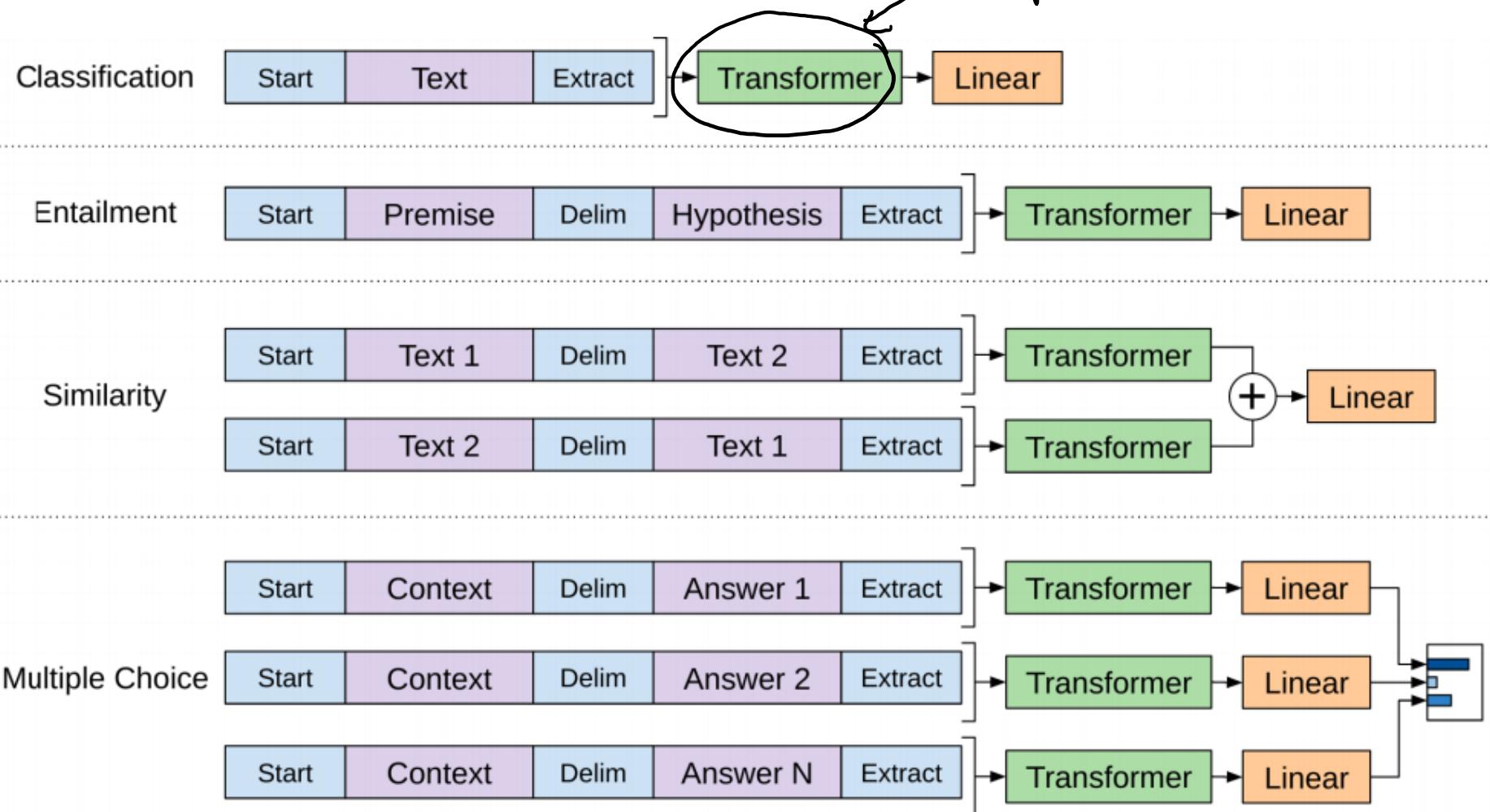


Table 2: Experimental results on natural language inference tasks, comparing our model with current state-of-the-art methods. 5x indicates an ensemble of 5 models. All datasets use accuracy as the evaluation metric.

Method	MNLI-m	MNLI-mm	SNLI	SciTail	QNLI	RTE
ESIM + ELMo [44] (5x)	-	-	<u>89.3</u>	-	-	-
CAFE [58] (5x)	80.2	79.0	<u>89.3</u>	-	-	-
Stochastic Answer Network [35] (3x)	<u>80.6</u>	<u>80.1</u>	-	-	-	-
CAFE [58]	78.7	77.9	88.5	<u>83.3</u>	-	-
GenSen [64]	71.4	71.3	-	-	<u>82.3</u>	59.2
Multi-task BiLSTM + Attn [64]	72.2	72.1	-	-	<u>82.1</u>	61.7
Finetuned Transformer LM (ours)	82.1	81.4	89.9	88.3	88.1	56.0

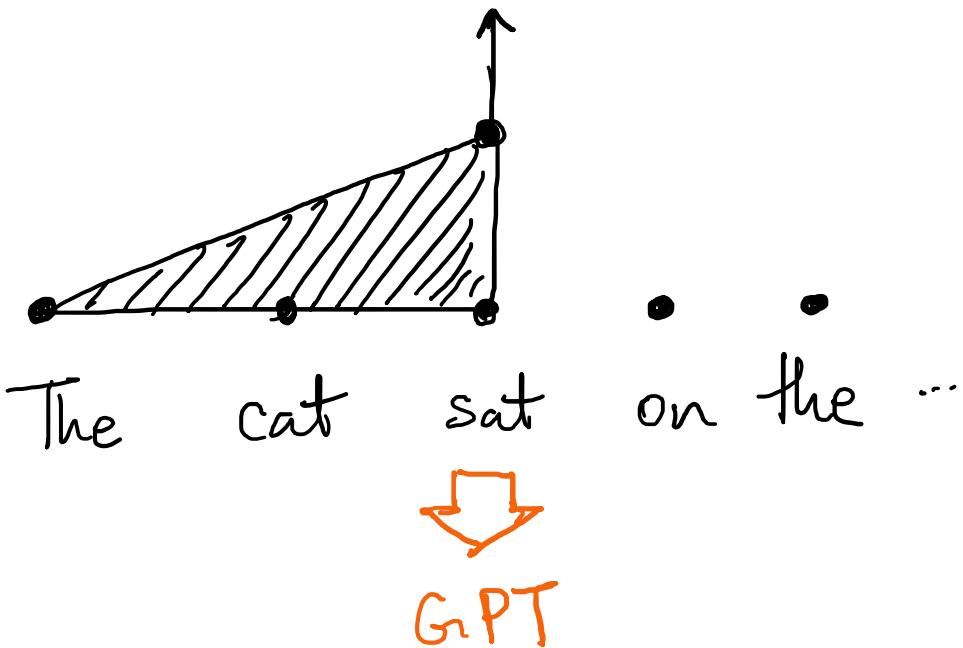
Table 3: Results on question answering and commonsense reasoning, comparing our model with current state-of-the-art methods.. 9x means an ensemble of 9 models.

Method	Story Cloze	RACE-m	RACE-h	RACE
val-LS-skip [55]	76.5	-	-	-
Hidden Coherence Model [7]	<u>77.6</u>	-	-	-
Dynamic Fusion Net [67] (9x)	-	55.6	49.4	51.2
BiAttention MRU [59] (9x)	-	60.2	<u>50.3</u>	<u>53.3</u>
Finetuned Transformer LM (ours)	86.5	62.9	57.4	59.0

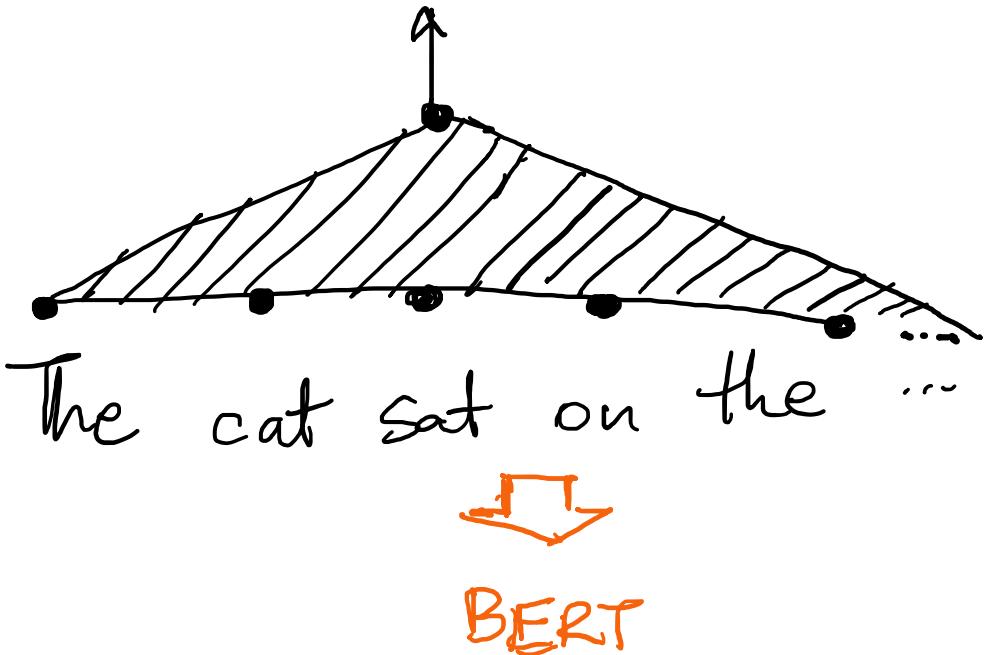
Results of
GPT-1
vs. the
then
state-of-the-
art.

Moving to BERT:

Self-attention

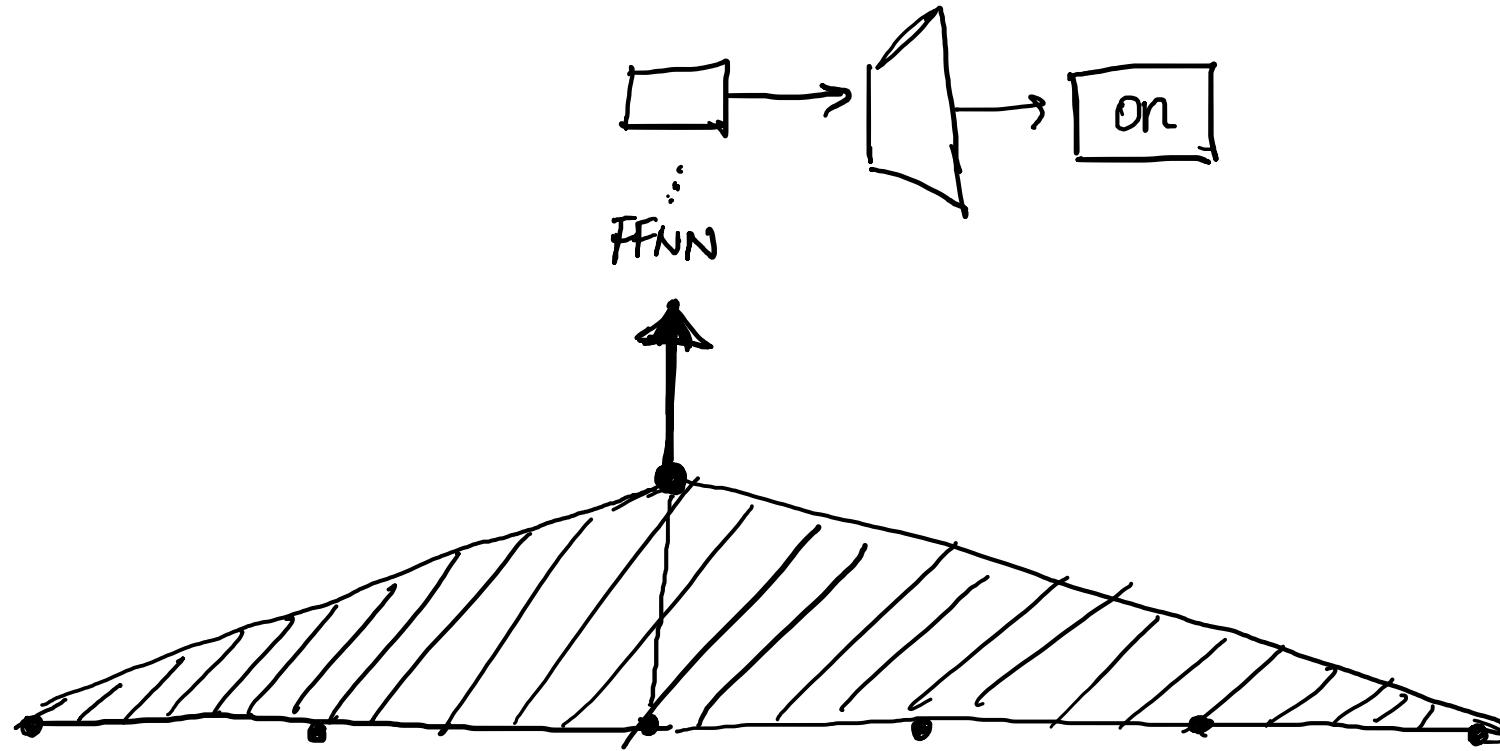


Masked self-attention



BERT:

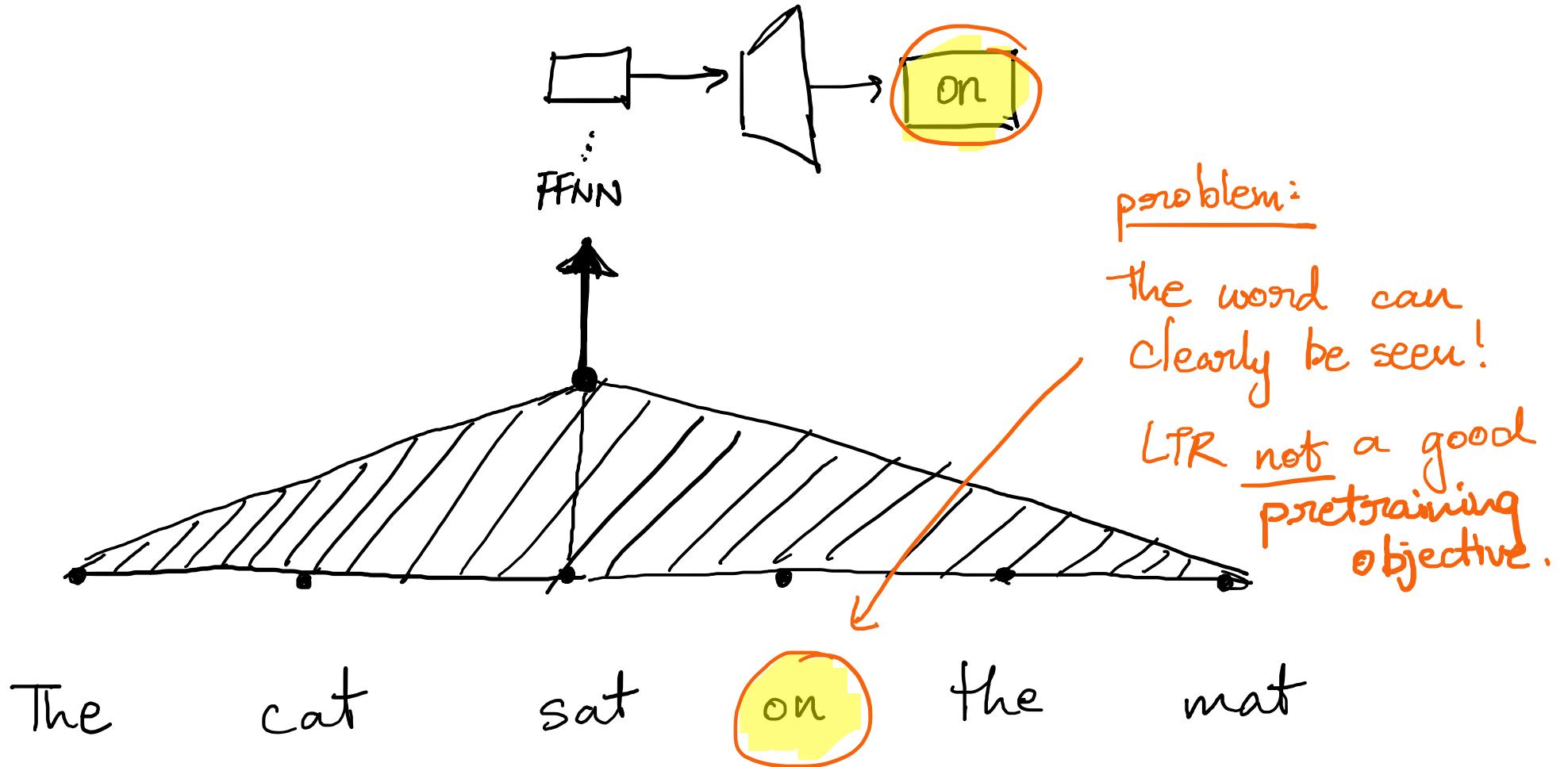
Stack encoders, not decoders!



The cat sat on the mat

BERT:

Stack encoders, not decoders!



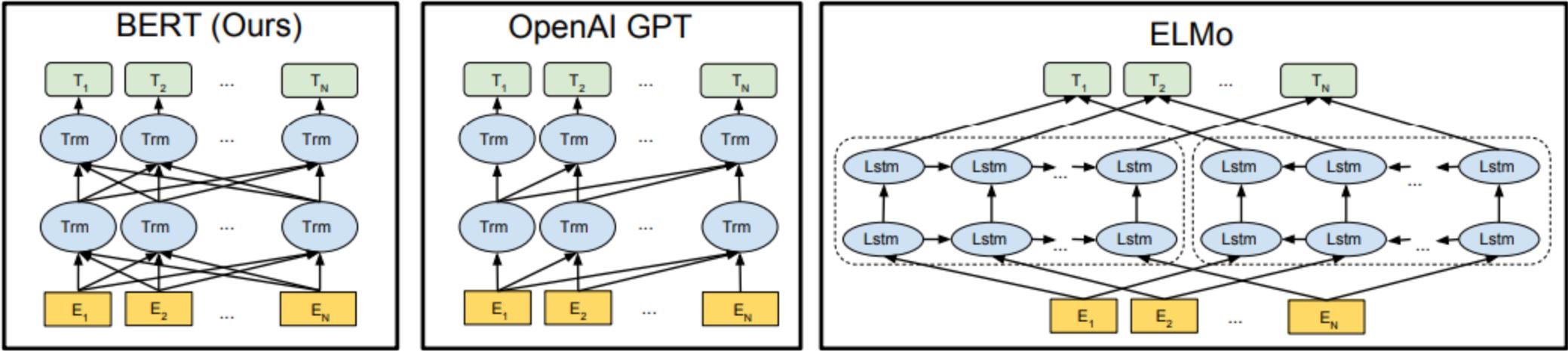


Figure 3: Differences in pre-training model architectures. BERT uses a bidirectional Transformer. OpenAI GPT uses a left-to-right Transformer. ELMo uses the concatenation of independently trained left-to-right and right-to-left LSTMs to generate features for downstream tasks. Among the three, only BERT representations are jointly conditioned on both left and right context in all layers. In addition to the architecture differences, BERT and OpenAI GPT are fine-tuning approaches, while ELMo is a feature-based approach.

BERT Pre-training objectives

① Masked language model (MLM)

Replace 15% words with [MASK] & predict them.

② Next sentence prediction (NSP).

Predict if sentence pair $\langle s_1, s_2 \rangle$ are a pair of adjacent sentences.

Masked Language Model.

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

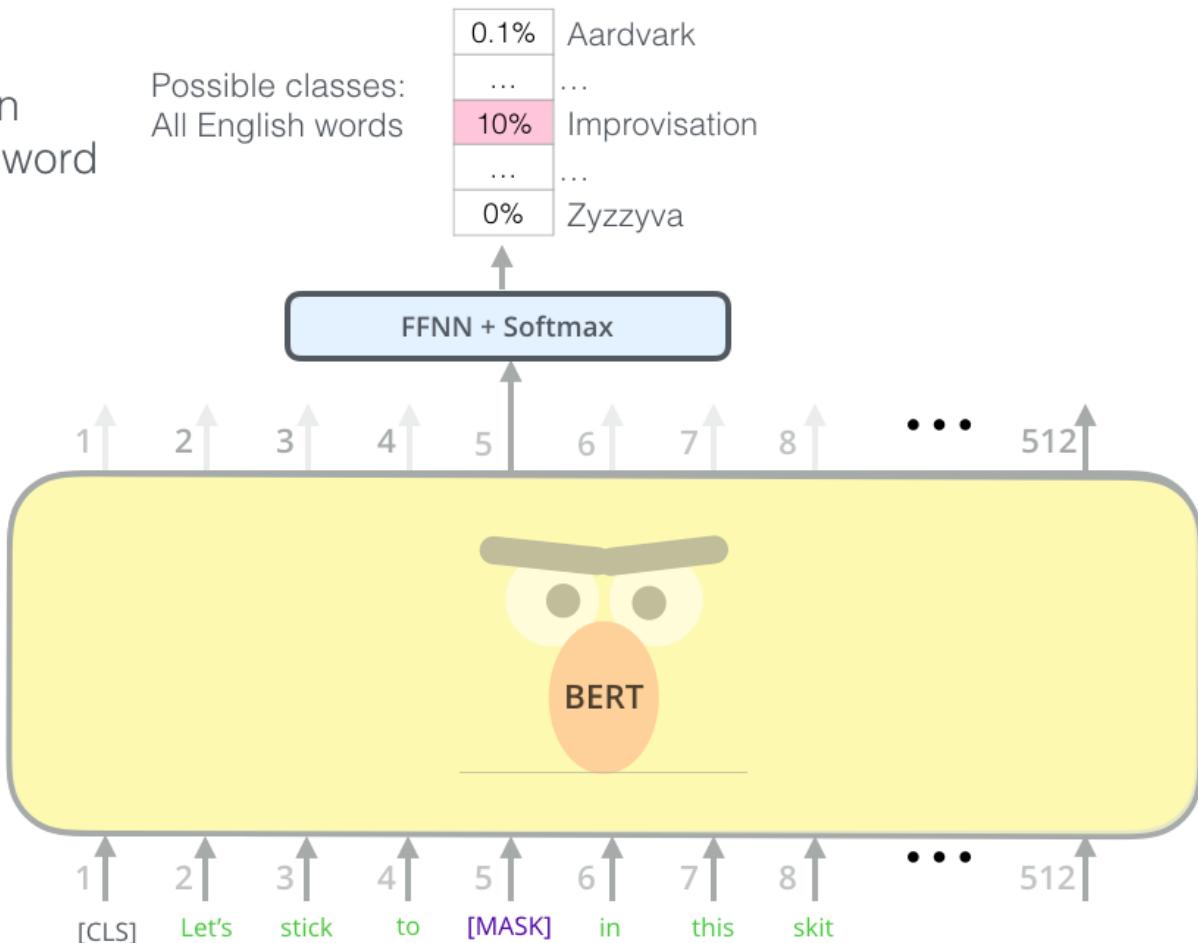
Possible classes:
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzyva

FFNN + Softmax

Randomly mask 15% of tokens

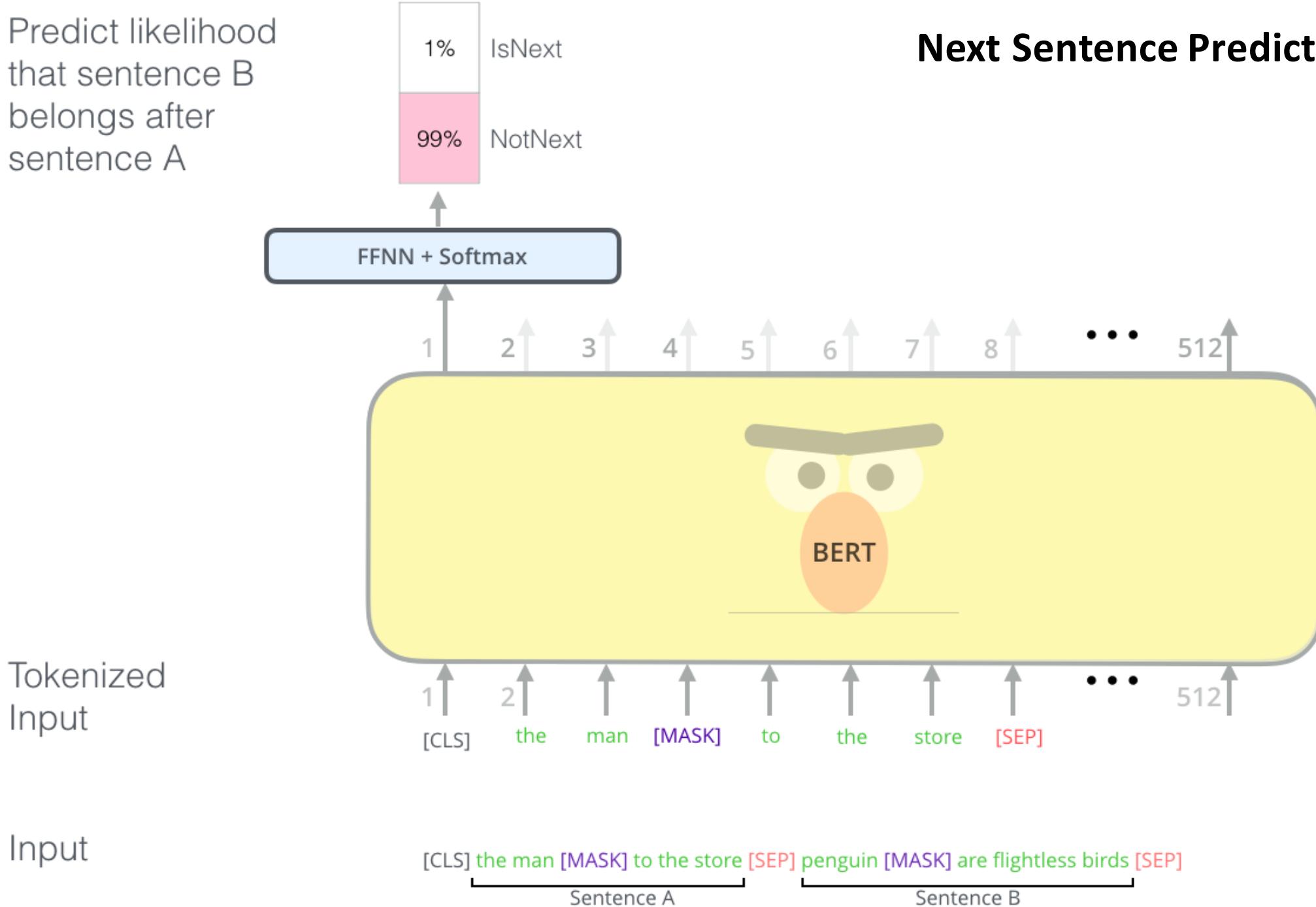
Input



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

Predict likelihood
that sentence B
belongs after
sentence A

Next Sentence Prediction (NSP)



Sentence Encodings

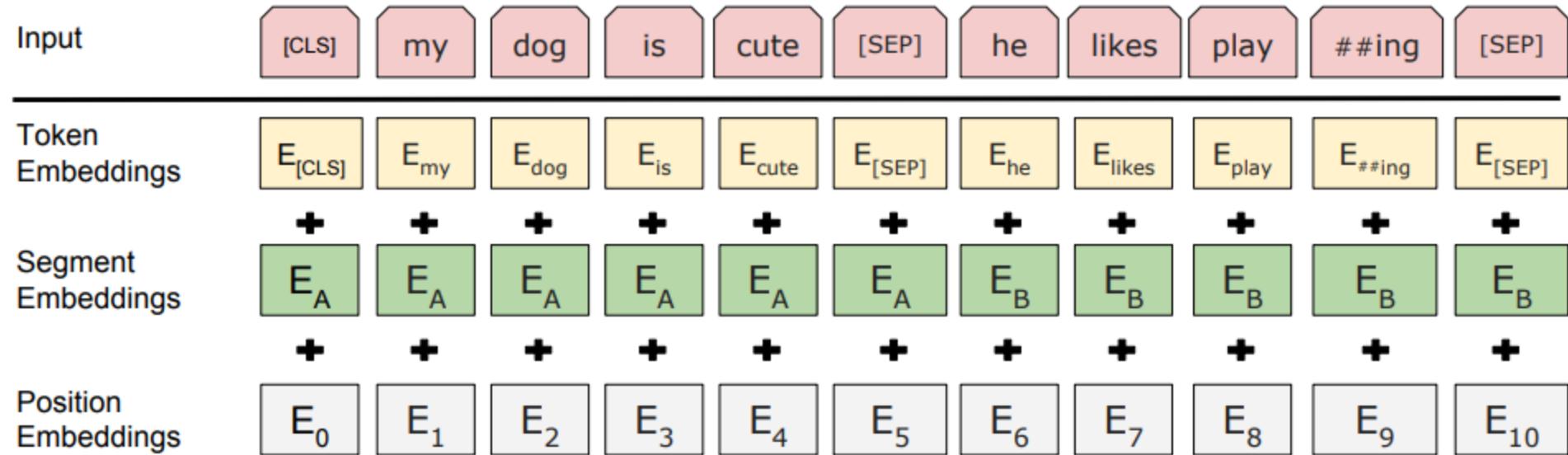
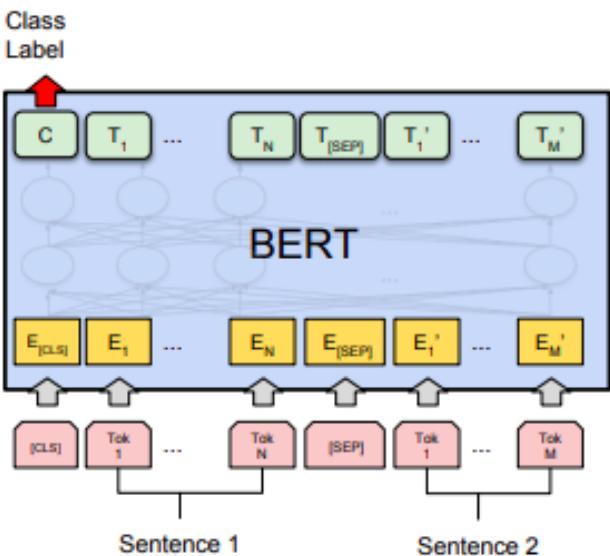
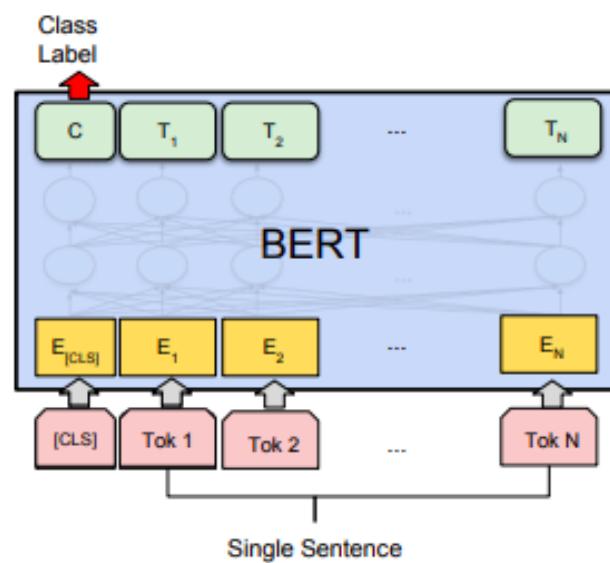


Figure 2: BERT input representation. The input embeddings are the sum of the token embeddings, the segmentation embeddings and the position embeddings.

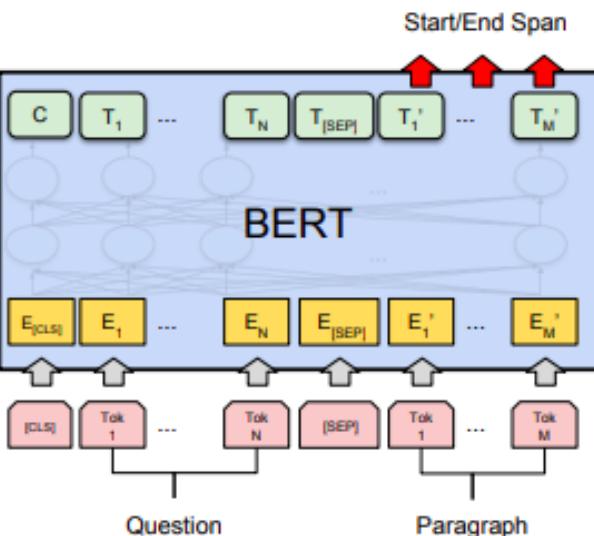
Task-specific Deployments



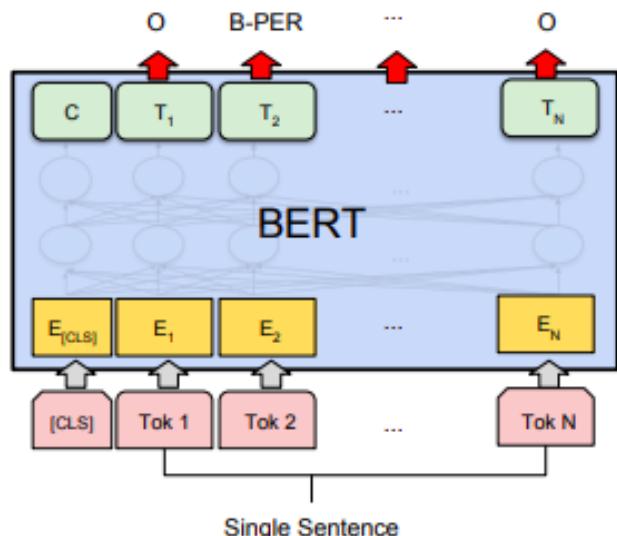
(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



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



(c) Question Answering Tasks:
SQuAD v1.1



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

System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1

Table 1: GLUE Test results, scored by the evaluation server (<https://gluebenchmark.com/leaderboard>). The number below each task denotes the number of training examples. The “Average” column is slightly different than the official GLUE score, since we exclude the problematic WNLI set.⁸ BERT and OpenAI GPT are single-model, single task. F1 scores are reported for QQP and MRPC, Spearman correlations are reported for STS-B, and accuracy scores are reported for the other tasks. We exclude entries that use BERT as one of their components.

System	Dev		Test	
	EM	F1	EM	F1
Top Leaderboard Systems (Dec 10th, 2018)				
Human	-	-	82.3	91.2
#1 Ensemble - nlnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	90.5
Published				
BiDAF+ELMo (Single)	-	85.6	-	85.8
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours				
BERT _{BASE} (Single)	80.8	88.5	-	-
BERT _{LARGE} (Single)	84.1	90.9	-	-
BERT _{LARGE} (Ensemble)	85.8	91.8	-	-
BERT _{LARGE} (Sgl.+TriviaQA)	84.2	91.1	85.1	91.8
BERT _{LARGE} (Ens.+TriviaQA)	86.2	92.2	87.4	93.2

Table 2: SQuAD 1.1 results. The BERT ensemble is 7x systems which use different pre-training checkpoints and fine-tuning seeds.

System	Dev		Test	
	EM	F1	EM	F1
Top Leaderboard Systems (Dec 10th, 2018)				
Human	86.3	89.0	86.9	89.5
#1 Single - MIR-MRC (F-Net)	-	-	74.8	78.0
#2 Single - nlnet	-	-	74.2	77.1
Published				
unet (Ensemble)	-	-	71.4	74.9
SLQA+ (Single)	-	-	71.4	74.4
Ours				
BERT _{LARGE} (Single)	78.7	81.9	80.0	83.1

Table 3: SQuAD 2.0 results. We exclude entries that use BERT as one of their components.

System	Dev	Test
ESIM+GloVe	51.9	52.7
ESIM+ELMo	59.1	59.2
OpenAI GPT	-	78.0
BERT _{BASE}	81.6	-
BERT _{LARGE}	86.6	86.3
Human (expert) [†]	-	85.0
Human (5 annotations) [†]	-	88.0

Table 4: SWAG Dev and Test accuracies. [†]Human performance is measured with 100 samples, as reported in the SWAG paper.

Ablation Study 1:

Without next sentence prediction,

Without bidirectional self-attention

Tasks	Dev Set				
	MNLI-m (Acc)	QNLI (Acc)	MRPC (Acc)	SST-2 (Acc)	SQuAD (F1)
BERT _{BASE}	84.4	88.4	86.7	92.7	88.5
No NSP	83.9	84.9	86.5	92.6	87.9
LTR & No NSP	82.1	84.3	77.5	92.1	77.8
+ BiLSTM	82.1	84.1	75.7	91.6	84.9

Table 5: Ablation over the pre-training tasks using the BERT_{BASE} architecture. “No NSP” is trained without the next sentence prediction task. “LTR & No NSP” is trained as a left-to-right LM without the next sentence prediction, like OpenAI GPT. “+ BiLSTM” adds a randomly initialized BiLSTM on top of the “LTR + No NSP” model during fine-tuning.

Ablation Study 1 (cot'd):

Without bidirectional self-attention

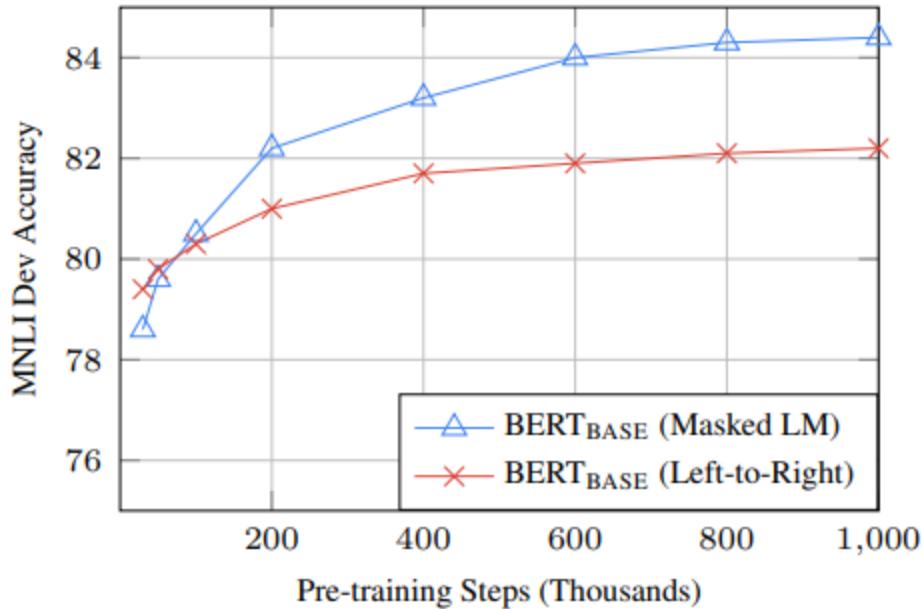


Figure 5: Ablation over number of training steps. This shows the MNLI accuracy after fine-tuning, starting from model parameters that have been pre-trained for k steps. The x-axis is the value of k .

Ablation Study 2: Transformer model hyperparameters

Hyperparams				Dev Set Accuracy		
#L	#H	#A	LM (ppl)	MNLI-m	MRPC	SST-2
3	768	12	5.84	77.9	79.8	88.4
6	768	3	5.24	80.6	82.2	90.7
6	768	12	4.68	81.9	84.8	91.3
12	768	12	3.99	84.4	86.7	92.9
12	1024	16	3.54	85.7	86.9	93.3
24	1024	16	3.23	86.6	87.8	93.7

Table 6: Ablation over BERT model size. #L = the number of layers; #H = hidden size; #A = number of attention heads. “LM (ppl)” is the masked LM perplexity of held-out training data.

Ablation Study 3: Masking strategy

Masking Rates			Dev Set Results		
MASK	SAME	RND	MNLI		NER
			Fine-tune	Fine-tune	Feature-based
80%	10%	10%	84.2	95.4	94.9
100%	0%	0%	84.3	94.9	94.0
80%	0%	20%	84.1	95.2	94.6
80%	20%	0%	84.4	95.2	94.7
0%	20%	80%	83.7	94.8	94.6
0%	0%	100%	83.6	94.9	94.6

Table 8: Ablation over different masking strategies.

[In the main paper]

Of the 15% words (tokens) to be masked:

1. 80% are changed to [mask]
2. 10% are kept unchanged
3. 10% are replaced with a random word

This is done to reduce discrepancy between pre-training and task-specific data.

[Here]

We vary these ratios to check what helps.

	MNLI	QNLI	QQP	RTE	SST	MRPC	CoLA	STS	WNLI	Avg
<i>Single-task single models on dev</i>										
BERT _{LARGE}	86.6/-	92.3	91.3	70.4	93.2	88.0	60.6	90.0	-	-
XLNet _{LARGE}	89.8/-	93.9	91.8	83.8	95.6	89.2	63.6	91.8	-	-
RoBERTa	90.2/90.2	94.7	92.2	86.6	96.4	90.9	68.0	92.4	91.3	-
<i>Ensembles on test (from leaderboard as of July 25, 2019)</i>										
ALICE	88.2/87.9	95.7	90.7	83.5	95.2	92.6	68.6	91.1	80.8	86.3
MT-DNN	87.9/87.4	96.0	89.9	86.3	96.5	92.7	68.4	91.1	89.0	87.6
XLNet	90.2/89.8	98.6	90.3	86.3	96.8	93.0	67.8	91.6	90.4	88.4
RoBERTa	90.8/90.2	98.9	90.2	88.2	96.7	92.3	67.8	92.2	89.0	88.5

Table 5: Results on GLUE. All results are based on a 24-layer architecture. BERT_{LARGE} and XLNet_{LARGE} results are from Devlin et al. (2019) and Yang et al. (2019), respectively. RoBERTa results on the development set are a median over five runs. RoBERTa results on the test set are ensembles of *single-task* models. For RTE, STS and MRPC we finetune starting from the MNLI model instead of the baseline pretrained model. Averages are obtained from the GLUE leaderboard.

Model		data	bsz	steps	SQuAD (v1.1/2.0)	MNLI-m	SST-2
RoBERTa							
with BOOKS + WIKI	16GB	8K	100K	93.6/87.3	89.0	95.3	
+ additional data (§3.2)	160GB	8K	100K	94.0/87.7	89.3	95.6	
+ pretrain longer	160GB	8K	300K	94.4/88.7	90.0	96.1	
+ pretrain even longer	160GB	8K	500K	94.6/89.4	90.2	96.4	
BERT_{LARGE}							
with BOOKS + WIKI	13GB	256	1M	90.9/81.8	86.6	93.7	
XLNet_{LARGE}							
with BOOKS + WIKI	13GB	256	1M	94.0/87.8	88.4	94.4	
+ additional data	126GB	2K	500K	94.5/88.8	89.8	95.6	

Table 4: Development set results for RoBERTa as we pretrain over more data (16GB → 160GB of text) and pretrain for longer (100K → 300K → 500K steps). Each row accumulates improvements from the rows above. RoBERTa matches the architecture and training objective of BERT_{LARGE}. Results for BERT_{LARGE} and XLNet_{LARGE} are from Devlin et al. (2019) and Yang et al. (2019), respectively. Complete results on all GLUE tasks can be found in the Appendix.

References:

Paper

<https://arxiv.org/pdf/1810.04805.pdf>

Google Colab Notebook

https://colab.research.google.com/github/tensorflow/tpu/blob/master/tools/colab/bert_finetuning_with_cloud_tpus.ipynb

GitHub repos

<https://github.com/google-research/bert> (TensorFlow 2.0)

<https://github.com/huggingface/pytorch-pretrained-BERT> (PyTorch)

Some illustrations from

<https://jalammar.github.io/illustrated-transformer/>

<https://jalammar.github.io/illustrated-gpt2/>

<https://jalammar.github.io/illustrated-bert/>

<https://www.topbots.com/generalized-language-models-bert-openai-gpt2/>