

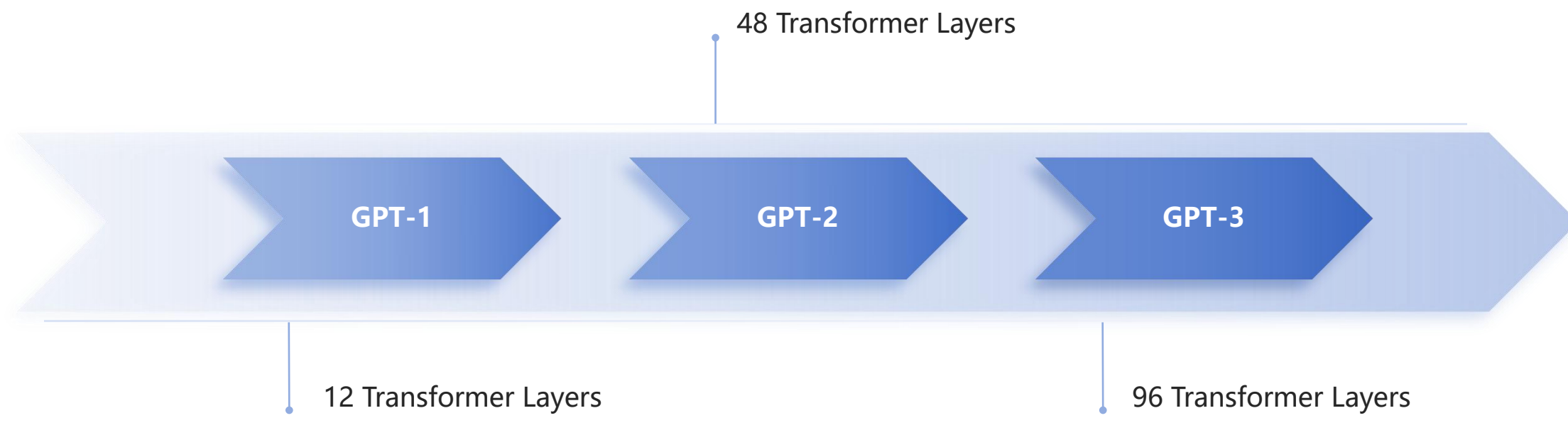


# From Static to Dynamic: A Deeper, Faster, and Adaptive Language Modeling Approach

Jiajia Li, Qiwei Li and Ping Wang\*

## Motivation

A key embodiment of the depth language model's reasoning and understanding ability for language is the depth of the model.



## Inference:

- Although the deeper the performance is better in overall performance, it is questionable whether each input instance requires the maximum number of layers.
- A short and simple sentence is obviously less difficult to encode than a complex or long sentence.
- It will risk overfitting if we use the same very deep structure for these sentences.
- The use of more structures than required for the encoding brings unnecessary computation overhead, thus increasing the service latency.

## Architecture

### (1) Vanilla Transformer Layer

For any  $i$ -th layer with input  $H_{i-1}$ , the encoding process for output  $H_i$  can be formalized as:

$$\tilde{H}_i = H_{i-1} + \text{MHATTN}(\text{LN}(H_{i-1}))$$

$$H_i = \tilde{H}_i + \text{FFN}(\text{LN}(\tilde{H}_i))$$

The computation of one-head in the multi-head attention of the vanilla Transformer can be represented as follows:

$$\text{AS}(X) = \frac{X^T \mathbf{W}_Q^T \mathbf{W}_K X}{\sqrt{d_{\text{head}}}}$$

$$\text{ATTN}(X) = \text{SOFTMAX}(\text{AS}(X)) \mathbf{W}_V X$$

Specifically, we calculate two additional positional scores:

$$\text{RKS}(X) = \mathbf{W}_K X \times \mathbf{E}_{R_{i-j}}, \text{RQS}(X) = \mathbf{W}_Q X \times \mathbf{E}_{R_{i-j}}$$

The new attention score:

$$\text{AS}^{\text{rel}}(X) = \text{AS}(X) + \text{RKS}(X) + \text{RQS}(X)$$

### (3) Structure Predictor

Perform first-token pooling on the outputs to obtain its sentence representation:

$$S = \text{POOLING}(H_{1Q})$$

Predict the number of layers in the hidden group it needs to use:

$$\text{probs} = \text{SOFTMAX}(\text{MLP}(S))$$

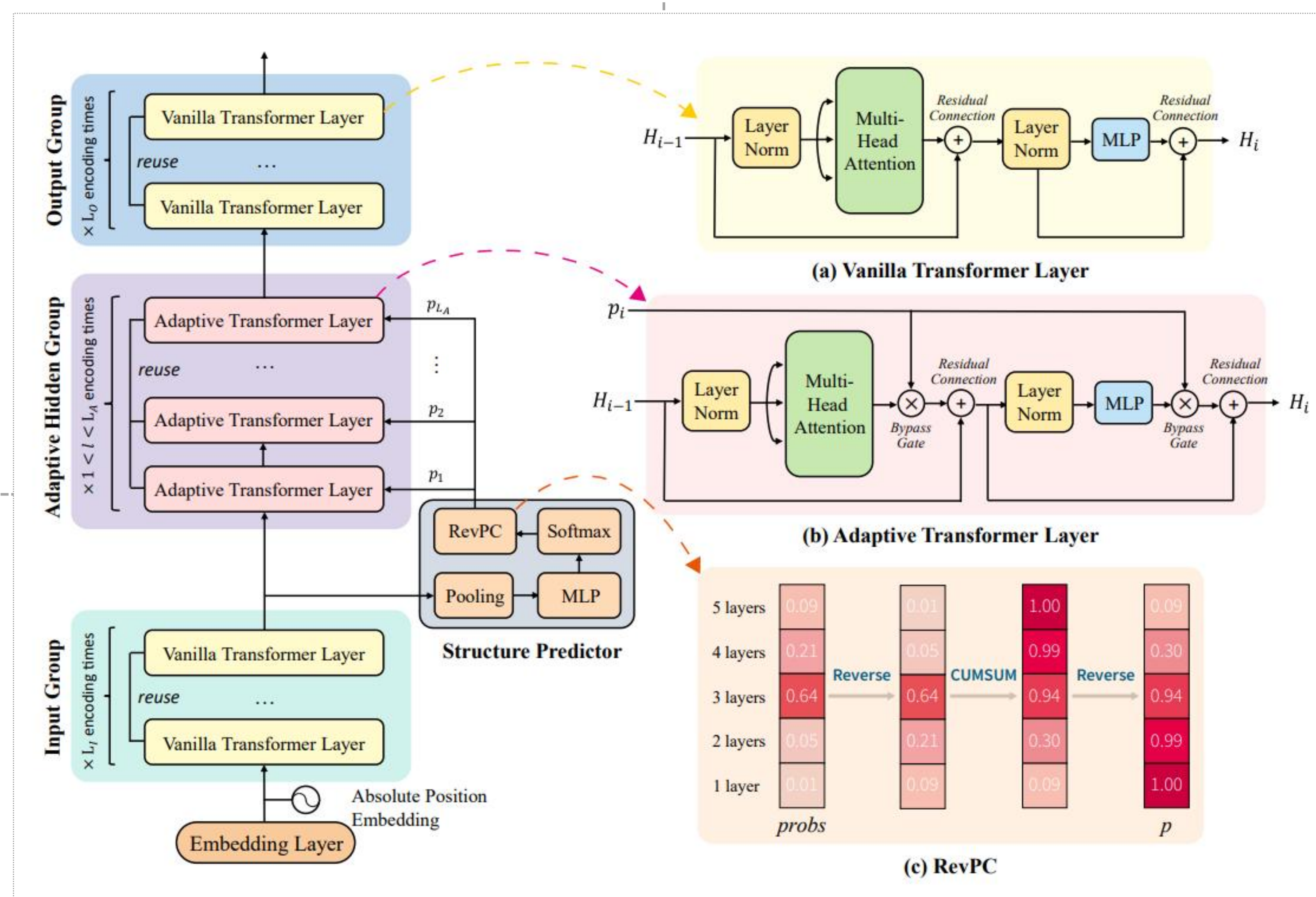
RevPC:

- 1) Reverse the probabilities.
- 2) Perform a cumulative sum of the elements.
- 3) Reverses the result back.

### (2) Adaptive Transformer Layer

Challenge of changing the model structure by directly changing the number of layers during the training phase:

- The operation of changing the model structure is not differentiable and cannot be directly optimized by the overall loss of the model.
- The training uses the mini-batch mechanism to take advantage of the parallel computing power of GPU devices and for the purpose of more stable optimization, however the model structure required for different examples within the mini-batch is not the same, which makes the input-dependent model structure unfeasible.



#### Advantage:

- The complexity estimation and model structure changes are optimized by the loss of the language model during pre-training, achieving self-learning of the input complexity without additional complexity learning process.
- The model can continue to update the complexity estimation according to the downstream finetuning in conjunction with the specific task, thus providing more flexibility

#### Process:

- During the training stage, it is possible to use mini-batch strategy to encode different examples simultaneously while each example actually uses a different model structure.
- In the inference stage, early exiting mechanism is directly performed according to the structure predicted by the complexity estimator on the input to achieve the purpose of efficient inference.

The encoding process of adaptive Transformer is changed from the vanilla Transformer:

$$\tilde{H}_i = H_{i-1} + \text{MHATTN}(\text{LN}(H_{i-1}))$$

$$H_i = \tilde{H}_i + \text{FFN}(\text{LN}(\tilde{H}_i))$$

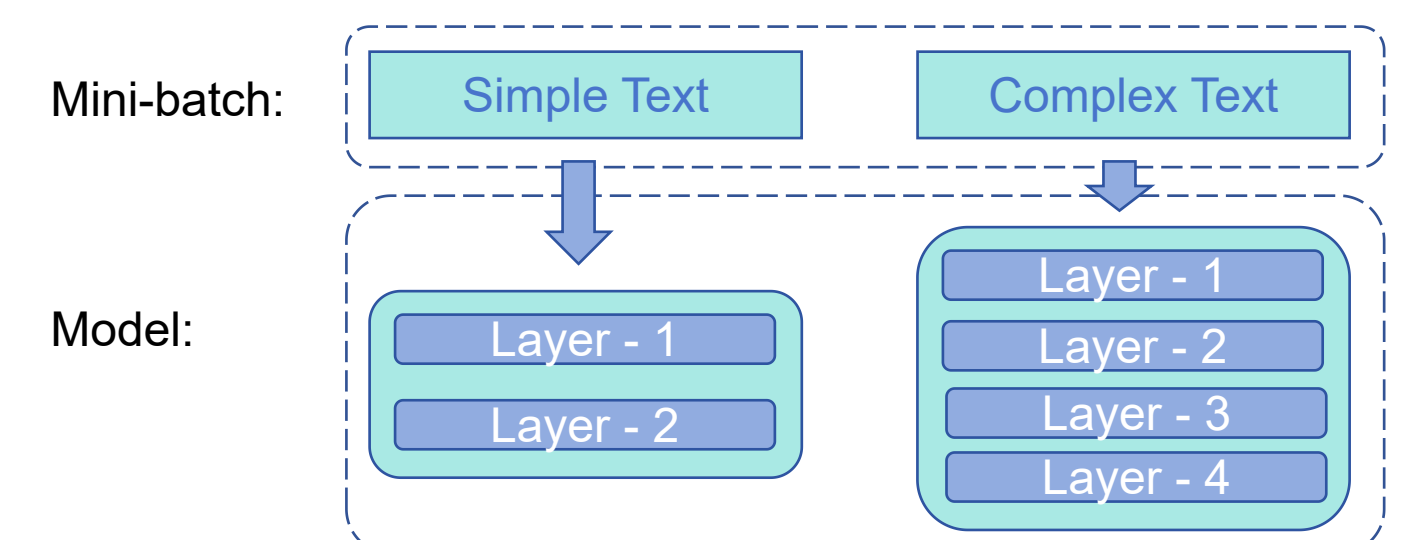
to:

$$\tilde{H}_i = H_{i-1} + p_i \cdot \text{MHATTN}(\text{LN}(H_{i-1}))$$

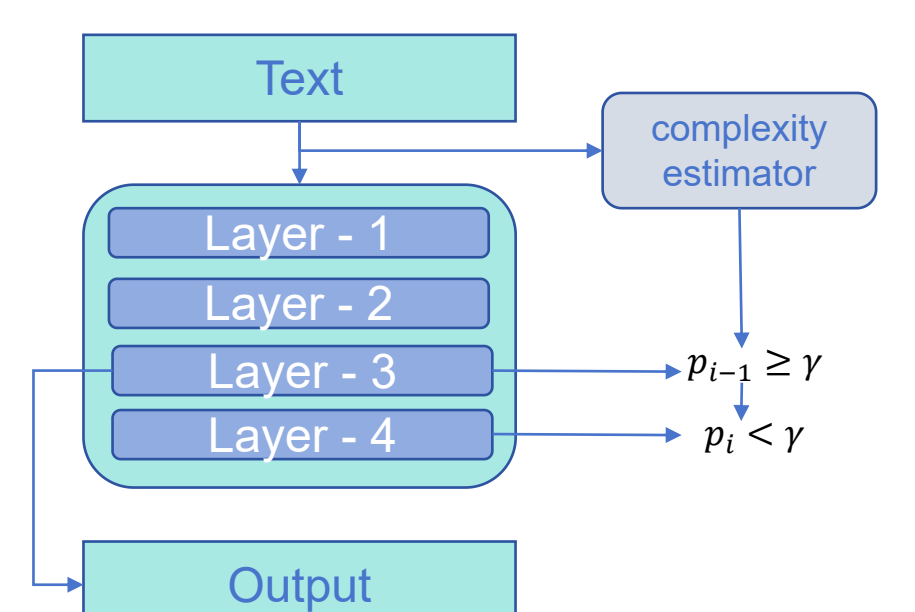
$$H_i = \tilde{H}_i + p_i \cdot \text{FFN}(\text{LN}(\tilde{H}_i))$$

### (4) Early Exiting

The training stage:



The inference stage:



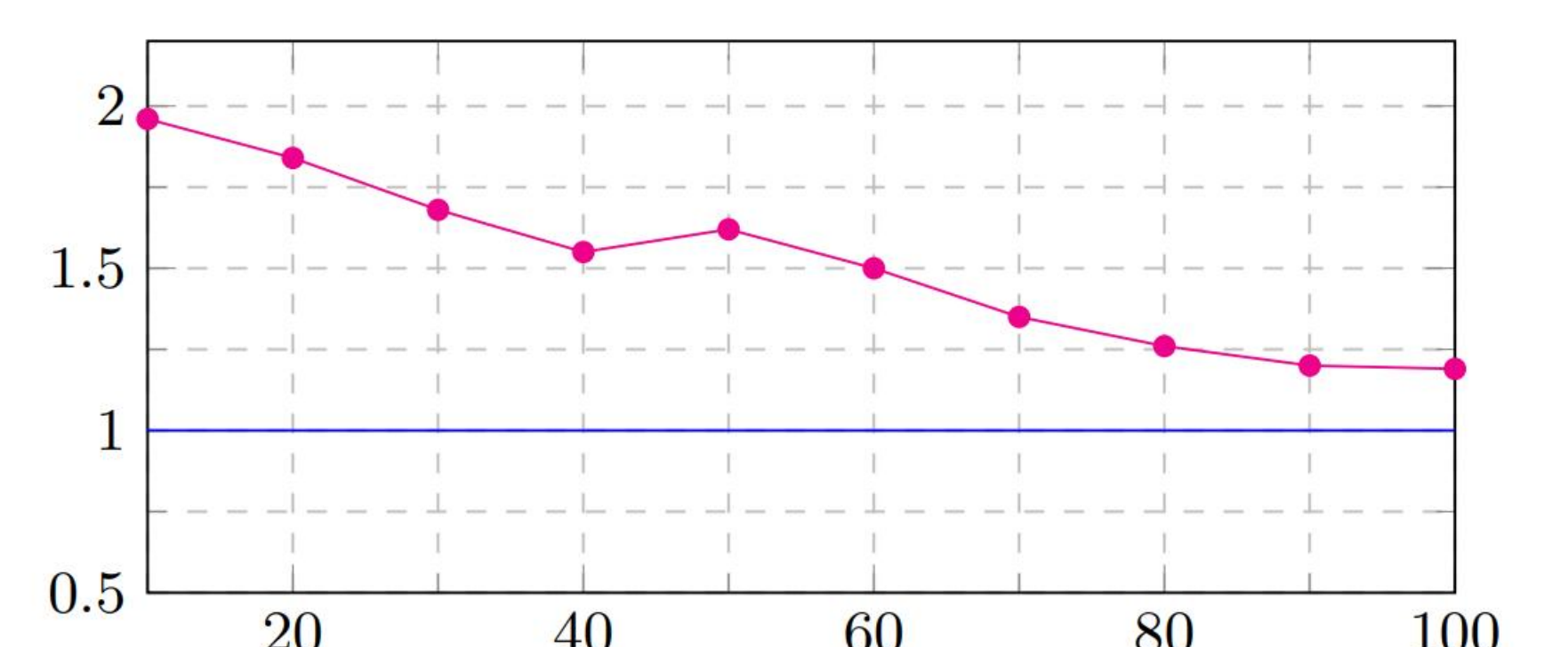
## Experiments

### Main Results

#### (1) Development Results on GLU NLU Benchmark

Model	Param.	MNLI	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Avg.	Speed
ALBERT-base	12M	84.6	89.6	91.7	92.8	58.9	89.5	89.5	78.6	84.4	(1.00×)
LayerDrop	12M	79.8	87.3	87.0	90.7	53.6	86.5	85.9	74.3	80.6	(1.96×)
HeadPrune	12M	80.3	88.0	86.8	90.5	54.1	87.4	86.2	75.1	81.1	(1.22×)
BranchyNet	12M	81.7	87.4	88.9	91.6	55.2	—	87.2	75.4	—	(1.88×)
Shallow-Deep	12M	81.5	87.8	89.2	91.7	55.5	—	87.1	75.2	—	(1.95×)
PABEE	12M	85.1	89.6	<b>91.8</b>	<b>93.0</b>	61.2	90.1	<b>90.0</b>	80.1	85.1	(1.57×)
<b>ALBERTa-base</b>	<b>30M</b>	<b>85.2</b>	<b>90.5</b>	91.6	<b>93.0</b>	<b>61.4</b>	<b>90.8</b>	89.9	<b>80.3</b>	<b>85.3</b>	(1.68×)

### The Consumed Time Statistics



Speedup vs. sequence lengths on sampled sentences.

#### (2) Development Results on SQUAD v2.0 MRC Benchmark

Model	Param.	EM	F1	Speedup
ALBERT-base	13M	77.1	80.0	(1.00×)
<b>ALBERTa-base</b>	<b>30M</b>	<b>79.3</b>	<b>82.6</b>	<b>(1.46×)</b>
ALBERT-large	19M	79.4	82.3	(1.00×)
<b>ALBERTa-large</b>	<b>48M</b>	<b>82.2</b>	<b>85.3</b>	<b>(1.50×)</b>
ALBERT-xlarge	62M	83.1	86.1	(1.00×)
<b>ALBERTa-xlarge</b>	<b>166M</b>	<b>84.2</b>	<b>87.3</b>	<b>(1.43×)</b>
ALBERT-xxlarge	237M	85.1	88.1	(1.00×)
<b>ALBERTa-xxlarge</b>	<b>630M</b>	<b>86.7</b>	<b>89.7</b>	<b>(1.55×)</b>

#### (3) Development Results on CoNLL-2003 NER Benchmark

Model	P	R	F1	Speed
ALBERT-base	93.69	94.10	93.90	(1.00×)
ALBERTa-base	94.15	94.27	<b>94.21</b>	<b>(1.81×)</b>
ALBERT-large	94.15	94.66	94.41	(1.00×)
ALBERTa-large	94.56	94.61	<b>94.59</b>	<b>(1.72×)</b>
ALBERT-xlarge	94.88	94.19	94.53	(1.00×)
ALBERTa-xlarge	94.53	95.22	<b>94.87</b>	<b>(1.70×)</b>
ALBERT-xxlarge	95.06	95.64	95.35	(1.00×)
ALBERTa-xxlarge	95.26	95.91	<b>95.58</b>	<b>(1.66×)</b>

### Ablation Study

Model	MNLI	QQP	QNLI	NER	Speed
<b>ALBERT-base</b>	84.6	89.6	91.7	93.90	(1.00×)
<b>ALBERTa-base</b>	85.2	90.5	91.6	94.21	(1.77×)
w/o EarlyExit	85.5	90.8	92.0	94.48	(0.28×)
$\gamma = 0.1$	85.2	90.7	92.0	94.37	(0.52×)
$\gamma = 0.5$	85.3	90.4	91.8	94.33	(1.22×)
$\gamma = 0.9$	80.4	87.5	88.2	91.23	(1.98×)
w/o RevPC	84.1	89.8	90.9	90.35	(1.71×)
w/o AbsPos	85.1	90.4	91.8	94.20	(1.78×)
w/o RelPos	84.9	90.0	91.5	94.05	(1.87×)