



# **RoBERTa: A Robustly Optimized BERT Pretraining Approach**

# **Idea: Improve BERT performance**

**B**idirectional **E**ncoder **R**epresentations from **T**ransformers

Popular self-training model

Goal: Increase performance with optimized pretraining



# BERT – Setup & Architecture

Takes two concatenated sequences of segments (more than one natural sentence)  $x_1, \dots, x_N$  and  $y_1, \dots, y_M$

$M + N < T$  (maximum sequence length)

Transformer architecture with layers  $L$ , self-attention heads  $A$  and hidden dimension  $H$



# BERT – Pretraining & Optimization

Two objectives:

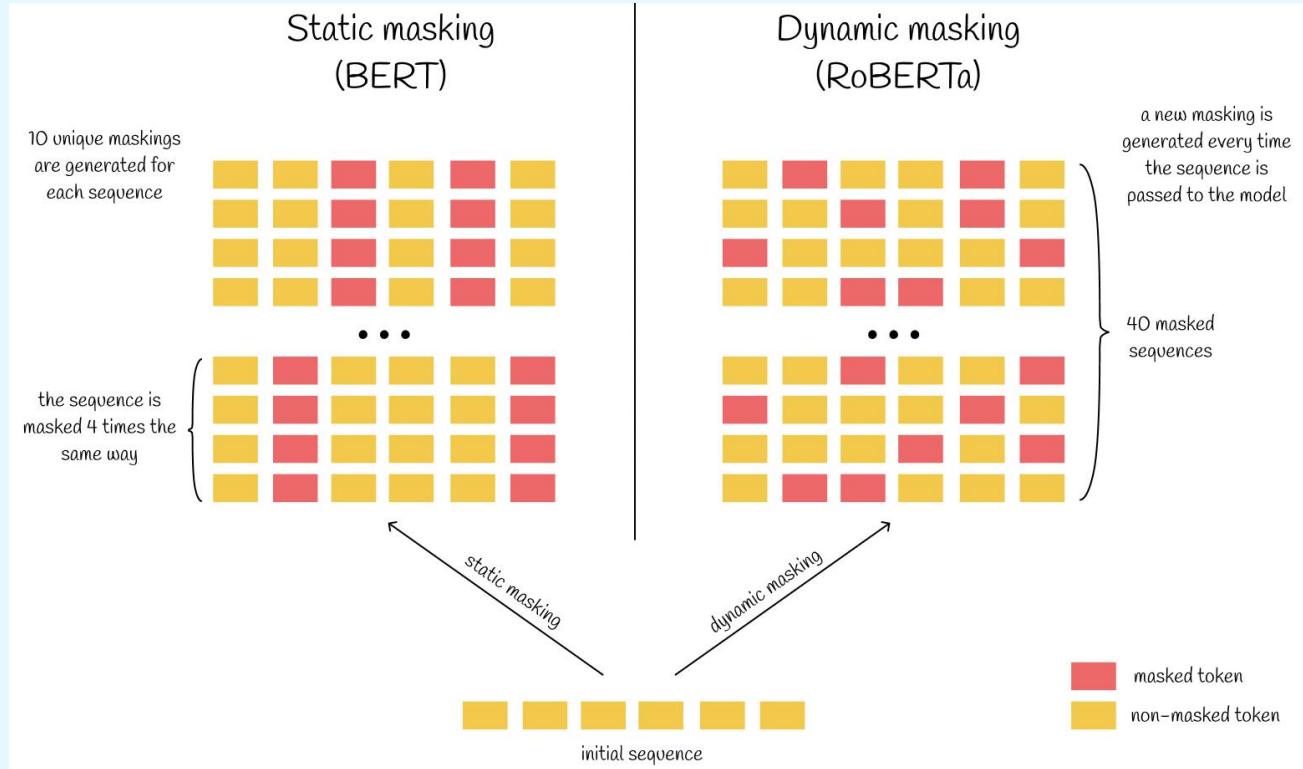
- Masked Language Model (MLM)
- Next Sentence Prediction (NSP)

Adam for optimization

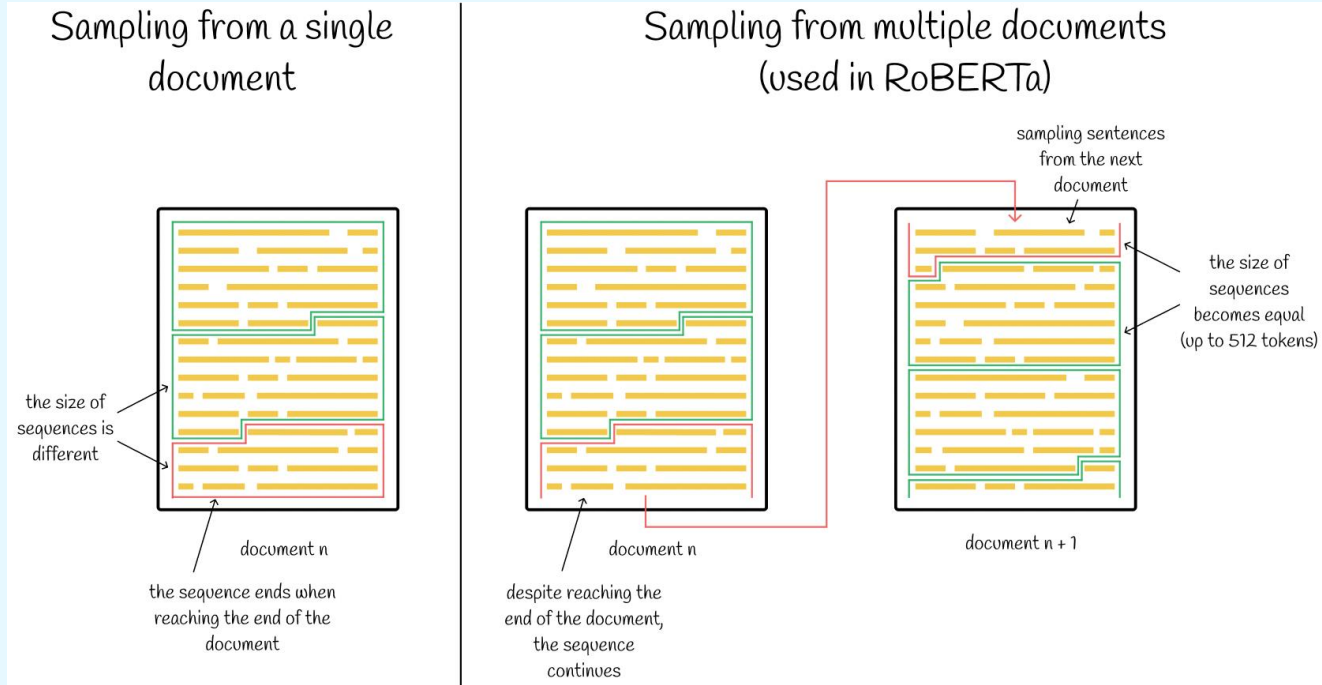
Trained with BookCorpus + English Wikipedia



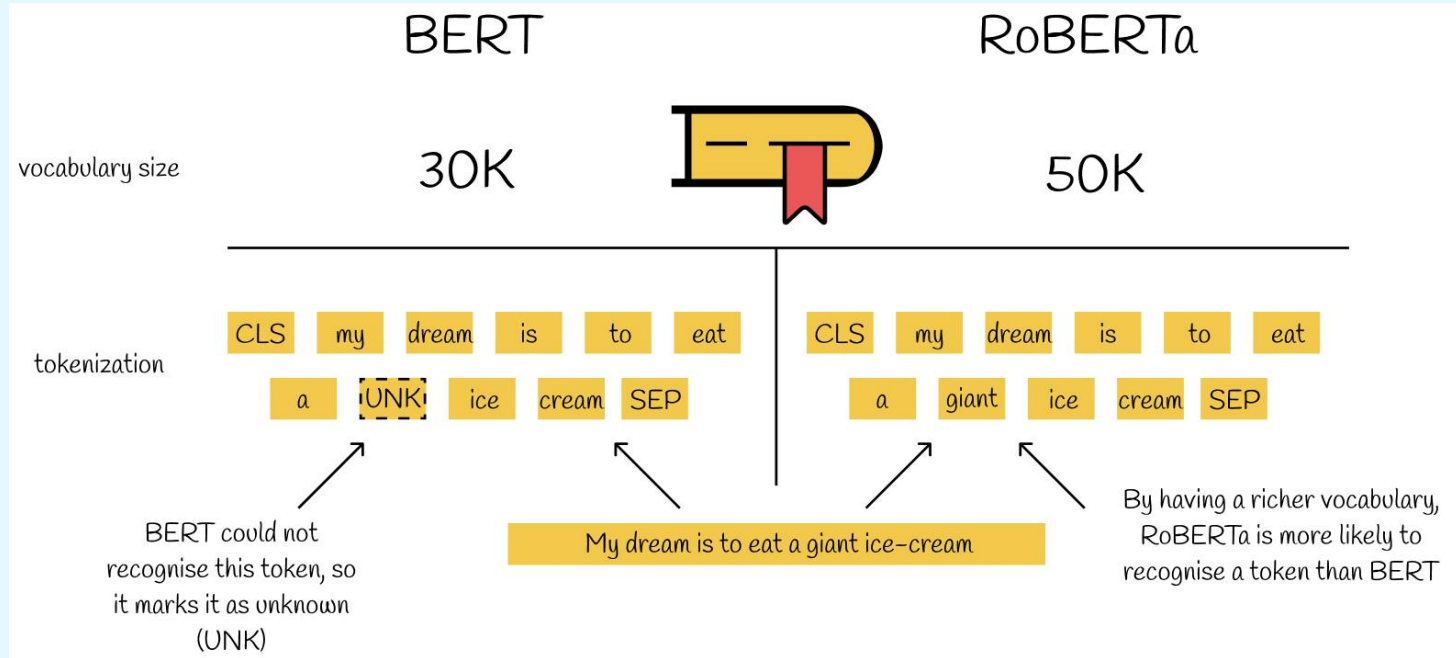
# RoBERTa improvements – Masking



# RoBERTa improvements – Training data



# RoBERTa improvements – Encoding



# RoBERTa improvements – Overview

- Dynamic masking
- Next sentence prediction
- Byte-pair encoding
- Larger batch size





# Evaluation

## Benchmarks

- GLUE
- SQuAD
- RACE

### General Language Understanding Evaluation

Contents	independent sentences and sentence pairs
Task	Answering questions, assigning correct context to ambiguous words ...
Challenge	Number of tasks

# Evaluation

## Benchmarks

- GLUE
- SQuAD
- RACE

### Stanford Question Answering Dataset

Contents	paragraph of context and a question
Task	Answer the question
Challenge	It contains unanswerable questions

# Evaluation

## Benchmarks

- GLUE
- SQuAD
- RACE

### ReAding Comprehension from Examinations

Contents	large dataset with 28,000 passages and nearly 100,000 questions
Task	select one correct answer from four options
Challenge	significantly longer context

# Evaluation

## Static vs. Dynamic Masking

Masking	SQuAD 2.0	MNLI-m	SST-2
reference	76.3	84.3	92.8
<i>Our reimplementation:</i>			
static	78.3	84.3	92.5
dynamic	78.7	84.0	92.9

# Evaluation

## Model Input Format and Next Sentence Prediction

Model	SQuAD 1.1/2.0	MNLI-m	SST-2	RACE
<i>Our reimplementation (with NSP loss):</i>				
SEGMENT-PAIR	90.4/78.7	84.0	92.9	64.2
SENTENCE-PAIR	88.7/76.2	82.9	92.1	63.0
<i>Our reimplementation (without NSP loss):</i>				
FULL-SENTENCES	90.4/79.1	84.7	92.5	64.8
DOC-SENTENCES	90.6/79.7	84.7	92.7	65.6
BERT <sub>BASE</sub>	88.5/76.3	84.3	92.8	64.3
XLNet <sub>BASE</sub> (K = 7)	-/81.3	85.8	92.7	66.1
XLNet <sub>BASE</sub> (K = 6)	-/81.0	85.6	93.4	66.7

# Evaluation

## Training with large batches

bsz	steps	lr	ppl	MNLI-m	SST-2
256	1M	1e-4	3.99	84.7	92.7
2K	125K	7e-4	<b>3.68</b>	<b>85.2</b>	<b>92.9</b>
8K	31K	1e-3	3.77	84.6	92.8

# Evaluation

## Text Encoding



# Differences to BERT

	BERT	RoBERTa
Masking	static masking	dynamic masking
Data training format	SEGMENT-PAIR	FULL-SENTENCES
Batch size (Pre training)	256	2048
vocab size	30K	50K
Data amount (Pre training)	3.3 billion word corpus	30.3 billion word corpus
Epochs (Pre training)	40	20
Steps (Pre training)	1,000,000	500,000



# BERT – RoBERTa comparison

	GLUE	SQuAD (v1.1/2.0)	RACE
BERT	84.05	90.9/81.8	72.0
RoBERTa	88.5	94.6/89.4	83.2

# Sources

Liu, Zhuang, et al. "A robustly optimized BERT pre-training approach with post-training." China National Conference on Chinese Computational Linguistics. Cham: Springer International Publishing, 2021.

Vyacheslav Efimov. "Large Language Models: RoBERTa — A Robustly Optimized BERT Approach" Toward Data Science, 25 September 2023, <https://towardsdatascience.com/roberta-1ef07226c8d8>. Accessed 9 January 2024.

# SQuAD

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	<b>94.6/89.4</b>	<b>90.2</b>	<b>96.4</b>
BERT <sub>LARGE</sub>						
with BOOKS + WIKI	13GB	256	1M	90.9/81.8	86.6	93.7
XLNet <sub>LARGE</sub>						
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

# GLUE

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

# SQuAD

Model	SQuAD 1.1		SQuAD 2.0	
	EM	F1	EM	F1
<i>Single models on dev, w/o data augmentation</i>				
BERT <sub>LARGE</sub>	84.1	90.9	79.0	81.8
XLNet <sub>LARGE</sub>	<b>89.0</b>	94.5	86.1	88.8
RoBERTa	88.9	<b>94.6</b>	<b>86.5</b>	<b>89.4</b>
<i>Single models on test (as of July 25, 2019)</i>				
XLNet <sub>LARGE</sub>			86.3 <sup>†</sup>	89.1 <sup>†</sup>
RoBERTa			86.8	89.8
XLNet + SG-Net Verifier			<b>87.0<sup>†</sup></b>	<b>89.9<sup>†</sup></b>

# RACE

Model	Accuracy	Middle	High
<i>Single models on test (as of July 25, 2019)</i>			
BERT <sub>LARGE</sub>	72.0	76.6	70.1
XLNet <sub>LARGE</sub>	81.7	85.4	80.2
RoBERTa	<b>83.2</b>	<b>86.5</b>	<b>81.3</b>