# RoBERTa: A Robustly Optimized BERT Pretraining Approach

**Yinhan Liu**[*§]    **Myle Ott**[*§]    **Naman Goyal**[*§]    **Jingfei Du**[*§]    **Mandar Joshi**[†]
**Danqi Chen**[§]    **Omer Levy**[§]    **Mike Lewis**[§]    **Luke Zettlemoyer**[†§]    **Veselin Stoyanov**[§]

[†] Paul G. Allen School of Computer Science & Engineering,
University of Washington, Seattle, WA
`{mandar90,lsz}@cs.washington.edu`

[§] Facebook AI
`{yinhanliu,myleott,naman,jingfeidu,`
`danqi,omerlevy,mikelewis,lsz,ves}@fb.com`

# Abstract

本篇論文在做的事情: 重新train 一個 優化版的Bert

相關資料: a replication study of BERT pretraining (Devlin et al. , 2019) that carefully measures the impact of many key hyperparameters and training data size

緣起: BERT was significantly undertrained, and can match or exceed the performance of every model published after it

重點及目標: hyperparameter choices have significant impact on the final results

評分標準: GLUE, RACE and SQuAD

結果: highlight the importance of previously overlooked design choices, and raise questions about the source of recently reported improvements
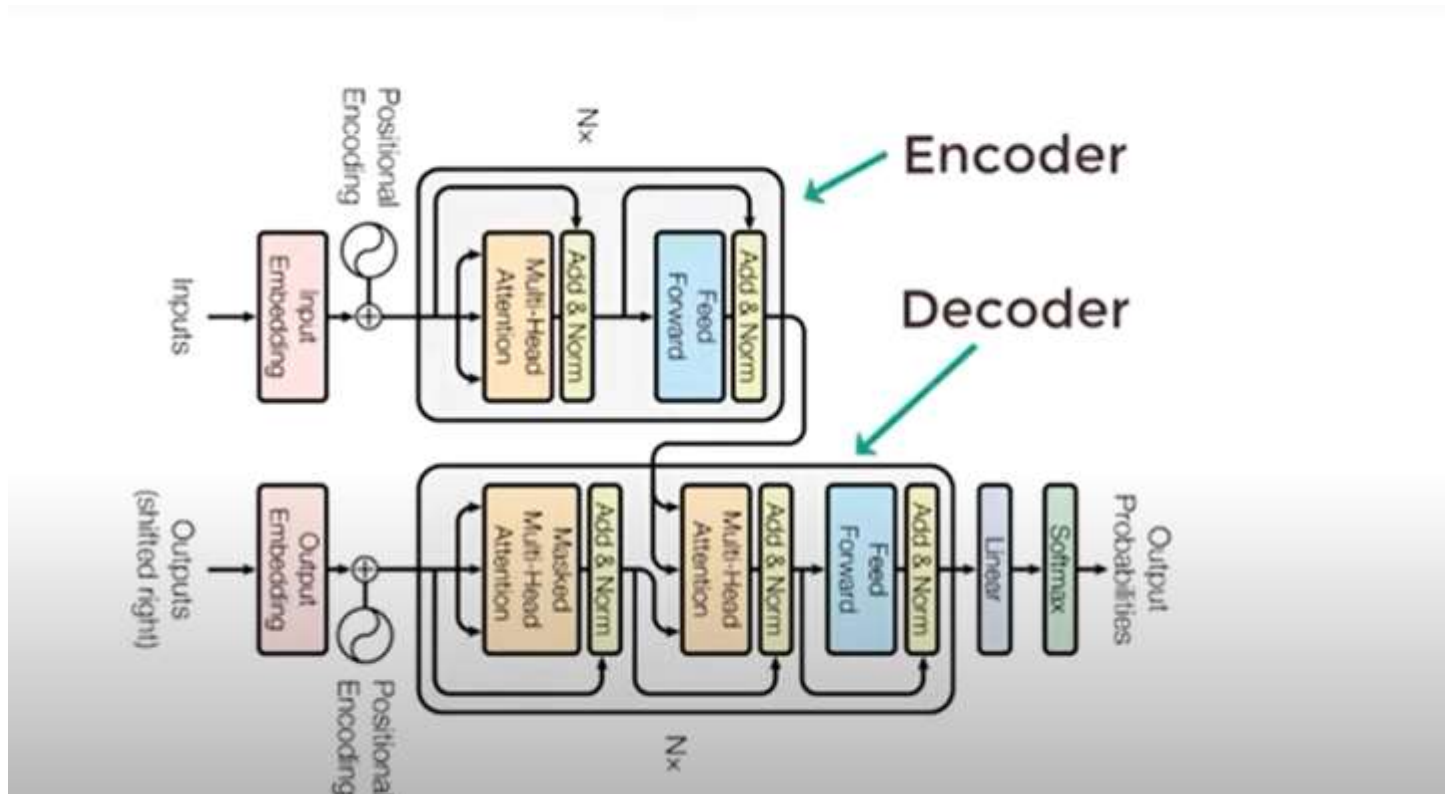
# Background

1. Architecture
2. Training Objectives
3. Optimization
4. Data

# Architecture



Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *North American Association for Computational Linguistics (NAACL)*.
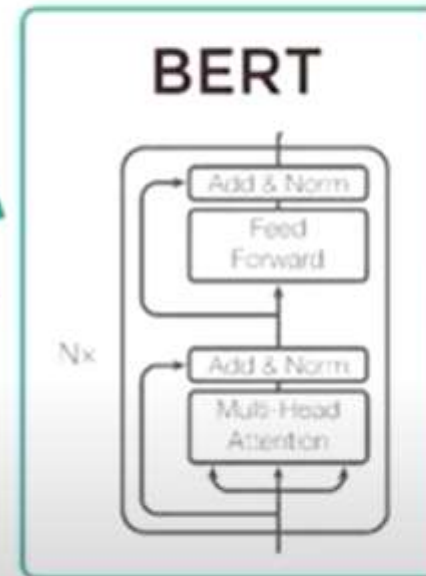
# Training Objectives

Pretraining (Pass 1) : "What is language? What is context?"

Masked Language Model (MLM)

The [MASK1] brown fox [MASK2] over the lazy dog.

BERT

Add & Norm
Feed Forward

Nx

Add & Norm
Multi-Head Attention

[MASK1] = quick
[MASK2]= jumped

Next Sentence Prediction (NSP)

A: Ajay is a cool dude.
B: He lives in Ohio

Yes. Sentence B follows sentence A

# Optimization

## Optimizer: Adam
## Activation Function: GEL

Dan Hendrycks and Kevin Gimpel. 2016. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415.*

# Data



BookCorpus (800M words) + Wikipedia (2500M words)

BOOKCORPUS (Zhu et al., 2015) plus English WIKIPEDIA, which totals 16GB of uncompressed text

# Tranining Procedure Analysis

1. static vs dynamic tasking
2. Input format and next sentence prediction
3. Training with large batches
4. Text Encoding (Byte-Pair Encoding(BPE) )

# static vs dynamic tasking

| Masking | SQuAD 2.0 | MNLI-m | SST-2 |
|---|---|---|---|
| reference | 76.3 | 84.3 | 92.8 |
| *Our reimplementation:* | | | |
| static | 78.3 | 84.3 | 92.5 |
| dynamic | 78.7 | 84.0 | 92.9 |

# Model Input Format and Next Sentence Prediction

NSP: Next Sentence Prediction Loss

1.SEGMENT-PAIR+NSP (用片段TRAIN)使用NSP LOSS

2.SENTENCE-PAIR+NSP (用句子TRAIN)使用NSP LOSS

3.FULL-SENTENCES不使用NSP LOSS

4.DOC-SENTENCES不使用NSP LOSS

| Model | SQuAD 1.1/2.0 | MNLI-m | SST-2 | RACE |
|---|---|---|---|---|
| *Our reimplementation (with NSP loss):* | | | | |
| SEGMENT-PAIR | 90.4/78.7 | 84.0 | 92.9 | 64.2 |
| SENTENCE-PAIR | 88.7/76.2 | 82.9 | 92.1 | 63.0 |
| *Our reimplementation (without NSP loss):* | | | | |
| 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$_{BASE}$ | 88.5/76.3 | 84.3 | 92.8 | 64.3 |
| XLNet$_{BASE}$ (K = 7) | –/81.3 | 85.8 | 92.7 | 66.1 |
| XLNet$_{BASE}$ (K = 6) | –/81.0 | 85.6 | 93.4 | 66.7 |

# 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 | **3.68** | **85.2** | **92.9** |
| 8K | 31K | 1e-3 | 3.77 | 84.6 | 92.8 |

# Text Encoding (Byte-Pair Encoding(BPE) )

Subword Tokenization: Byte Pair Encoding – YouTube

[1909.03341] Neural Machine Translation with Byte-Level Subwords (arxiv.org)

# RoBERTa

# Overall Result

| 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<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

General Language Understanding Evaluation (GLUE)

https://gluebenchmark.com/

- **Corpus of Linguistic Acceptability (CoLA)**
- **Stanford Sentiment Treebank (SST-2)**
- **Microsoft Research Paraphrase Corpus (MRPC)**
- **Quora Question Pairs (QQP)**
- **Semantic Textual Similarity Benchmark (STS-B)**
- **Multi-Genre Natural Language Inference (MNLI)**
- **Question-answering NLI (QNLI)**
- **Recognizing Textual Entailment (RTE)**
- **Winograd NLI (WNLI)**

GLUE also has Chinese version (https://www.cluebenchmarks.com/)

# Glue Result

| | MNLI | QNLI | QQP | RTE | SST | MRPC | CoLA | STS | WNLI | Avg |
|---|---|---|---|---|---|---|---|---|---|---|
| *Single-task single models on dev* | | | | | | | | | | |
| BERT$_{\text{LARGE}}$ | 86.6/- | 92.3 | 91.3 | 70.4 | 93.2 | 88.0 | 60.6 | 90.0 | - | - |
| XLNet$_{\text{LARGE}}$ | 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** | - |
| *Ensembles on test (from leaderboard as of July 25, 2019)* | | | | | | | | | | |
| 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** |

# SQuAD Introduction and Result

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

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity. The main forms of precipitation include drizzle, rain, sleet, snow, graupel and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".

What causes precipitation to fall?
gravity

What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?
graupel

Where do water droplets collide with ice crystals to form precipitation?
within a cloud

Sample question and answer pairs for a passage from the SQuAD dataset. Image credits to Rajpurkar et al., the original creators of the dataset.

# RACE Introduction and Result

ReAding Comprehension Dataset From Examinations

| Model | Accuracy | Middle | High |
|-------|----------|--------|------|
| *Single models on test (as of July 25, 2019)* | | | |
| BERT$_{LARGE}$ | 72.0 | 76.6 | 70.1 |
| XLNet$_{LARGE}$ | 81.7 | 85.4 | 80.2 |
| **RoBERTa** | **83.2** | **86.5** | **81.3** |

Passage: Do you love holidays but hate gaining weight? You are not alone. Holidays are times for celebrating. Many people are worried about their weight. With proper planning, though, it is possible to keep normal weight during the holidays. The idea is to enjoy the holidays but not to eat too much. You don't have to turn away from the foods that you enjoy.

Here are some tips for preventing weight gain and maintaining physical fitness:

Don't skip meals. Before you leave home, have a small, low-fat meal or snack. This may help to avoid getting too excited before delicious foods.

Control the amount of food. Use a small plate that may encourage you to "load up". You should be most comfortable eating an amount of food about the size of your fist.

Begin with soup and fruit or vegetables. Fill up beforehand on water-based soup and raw fruit or vegetables, or drink a large glass of water before you eat to help you to feel full.

Avoid high-fat foods. Dishes that look oily or creamy may have large amount of fat. Choose lean meat . Fill your plate with salad and green vegetables. Use lemon juice instead of creamy food.

Stick to physical activity. Don't let exercise take a break during the holidays. A 20-minute walk helps to burn off extra calories.

Questions:
What is the best title of the passage?
Options:
A. How to avoid holiday feasting
B. Do's and don'ts for keeping slim and fit.
C. How to avoid weight gain over holidays.
D. Wonderful holidays, boring experiences.