

Explore The Power of Transfer Learning: BERT for Toxic Comments Classification

Haoran Wang

University of Oregon

hwang8@cs.uoregon.edu

Abstract

BERT is a language model that was created and published in 2018 by Jacob Devlin and Ming-Wei Chang from Google [4]. It has taken the NLP landscape by a storm. BERT has been able to achieve state-of-the-art results in a wide variety of NLP tasks. Many researchers have been working on improving BERT in the past two years. As of now, we have many variations of the original BERT, such as ALBERT, RoBERTa, DistillBERT, and XLM-RoBERTa, etc. Transfer learning on those models has achieved state-of-the-art results when applied to different NLP tasks. In this paper, I show that transfer learning utilizing BERT and XLM-RoBERTa achieves near-top quality results on text classification in comparison to existing, simpler models such as CNN. Fine-tuned XLM-RoBERTa (cross lingual model), achieves an accuracy 0.9459 on the Jigsaw test set. In comparison, the top entry on Kaggle for Jigsaw scores 0.9555. Training was accelerated using Google's hardware TPU (tensor processing units) in conjunction with PyTorch-XLA. This results in an 4x speed increase over GPU training.

1 Introduction

For this project, I competed in the Jigsaw Multi-lingual Toxic Comment Classification competition hosted on Kaggle by Jigsaw and Google [6]. Jigsaw [5], formerly known as Google Ideas, is a unit within Google that tackles emerging online threats to make the internet a better place for all. With the increasing prevalence of social media, cyber bullying has become a real threat, especially to teenagers. Jigsaw and Google have therefore dedicated \$50,000 to host this Kaggle competition to explore how different NLP techniques can be used to classify toxic comments.

Text classification is one of the principal tasks

encountered in NLP. I applied transfer learning on pre-trained BERT and XLM-RoBERTa, and compare the results against CNN, which uses BERT word embeddings fed into a CNN.

In addition to fine-tuning BERT models, I used TPUs (tensor processing units) for training acceleration. TPUs are developed by Google, and only support Tensorflow natively. To enable PyTorch on TPU, I utilized the PyTorch/XLA package [10].

2 Background

This section will review the architecture of BERT and XLM-RoBERTa, as well as the concept of Transfer Learning.

2.1 BERT

BERT, short for **B**idirectional **E**ncoder **R**epresentation from **T**ransformers has two pre-trained models: BERT base and BERT large. BERT base has 12 layers, and around 110 million parameters. BERT large has 24 layers, and around 340 million parameters [4]. This is quite large. Therefore, I used BERT base for this project due to my limited computational resources.

To understand BERT, we need to first understand how Transformers work. BERT is basically layers of bidirectional Transformer encoders.

Transformers consist of a number of encoders and decoders. Each encoder and decoder has two layers: a multi-headed attention layer and a fully connected feed-forward neural network.

Multi-headed attention computes self-attention multiple times in parallel and independently. The outputs are then concatenated and linearly transformed.

“Self-attention, sometimes called intra-attention, is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. [11]” Therefore, it helps the encoder to look at other words in the input sentence when encoding each word. Thus, it makes the encoder take the contexts of the sentence into consideration.

To calculate self-attention, the encoder creates a query vector, a key vector, and a value vector for the embeddings of different words in the sequence. Next, it calculates the dot product of the query vector with the key vector for each word of the input sequence against the other words in the sequence. Then, it divides the dot product by the square root of the dimension of the key vectors, followed by an application of the softmax function. Finally it multiplies each value vector by the softmax score and sums up the weighted value vectors.

Transformers also add positional encoding when encoding each word. BERT’s structure makes BERT different from context-free word embeddings such as word2vec and GloVe, which simply generate word embedding (vectors) for each word regardless of that word’s context. Instead, BERT generates different word embeddings based on word context. This makes BERT a very powerful language model.

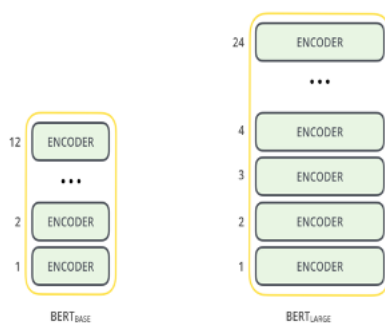


Figure 1: BERT base and BERT large [4]

2.2 Transfer Learning

Because of how well the BERT models performs, BERT is often combined with transfer learning to solve different NLP tasks. Instead of building a model from scratch, we can use a pre-trained

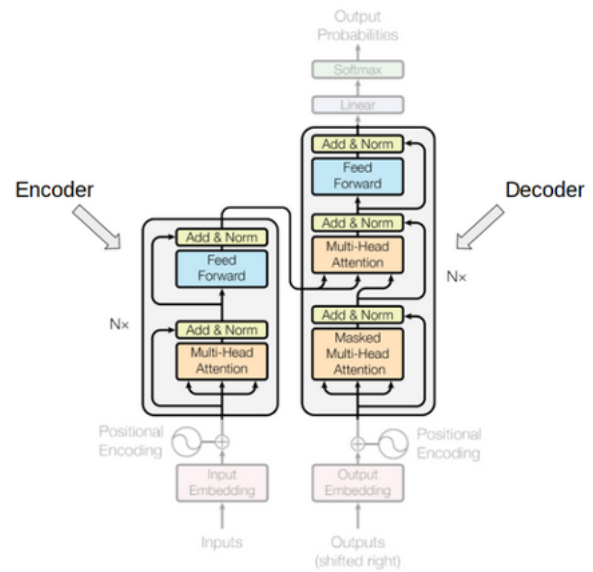


Figure 2: Transformer Architecture [11]

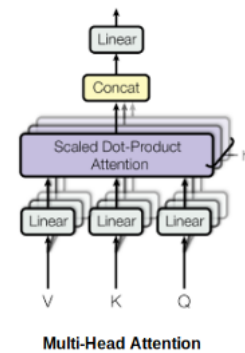


Figure 3: Multi-Head Attention [11]

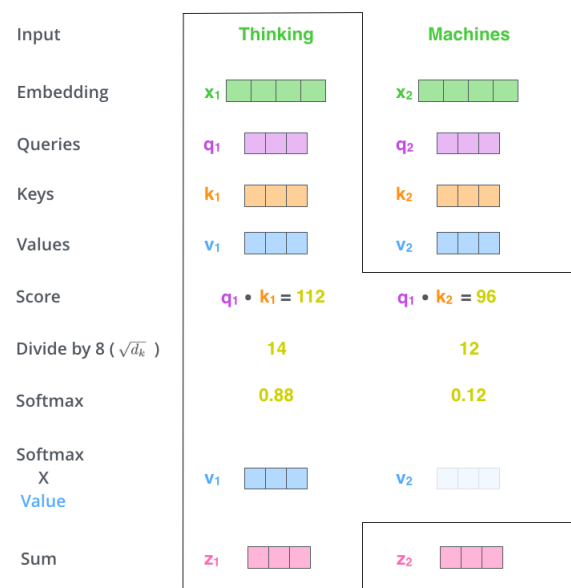


Figure 4: Calculate Self-Attention [1]

model like BERT to solve our own specific problems, with a little bit fine-tuning to adapt the model. This saves a lot of resources because our model does not need to learn every parameter from scratch. In the case of pre-trained BERT (base), we do not need to learn the first 12 layers of parameters. We can just add a layer on top of them to solve specific problems.

BERT is pre-trained on two tasks: Masked Language Modeling (MLM) and Next Sentence Prediction (NSP). For MLM, before feeding text sequences into BERT, some words in the sequence are masked with the [MASK] token. BERT then tries to predict the original value of the masked words based on the context of the sequence.

For NSP, BERT receives pairs of text sequences as input, and BERT tries to predict if the second sentence is the subsequent sentence in the original document. When the BERT model was trained, both MSM and NSP were trained together, and the goal was to minimize the combined loss function of both.

2.3 XLM-RoBERTa

BERT is trained entirely on English data. Many researchers who are working on non-English data have then come up with language specific models such as Finnish BERT, French BERT and German BERT. Research on a multilingual model continues to this day. In November 2019, the Facebook AI team released XLM-RoBERT [3], an update to their XLM-100 model. They used the RoBERTa [8] model (A Robustly Optimized BERT Pretraining Approach) to improve XLM-100 and hence produce XLM-RoBERTa.

XLM-RoBERTa is trained on 100 languages, total of 2.5TB of text. The pre-trained XLM-RoBERTa has 550 million parameters, compared to RoBERTa's 355 million.

2.4 TPU

The Tensor Processing Unit is an ASIC developed by Google to specialize in deep learning. Models that previously took days to train on GPU can be trained in hours on TPU. However, in order to use PyTorch on TPU, we have to use a package called Pytorch/XLA. There are currently no stable

releases of this platform.

3 Data

The Jigsaw Kaggle competition has provided a training set, validation set, and test set. The training set's comments are entirely in English. The test set's comments are composed of multiple non-English languages. I have combined the training set from the previous competition to make a large training set of 2.1 million comments. Since my computational resources are limited, I only used the first 200,000 datapoints as my training set.

I used the BERT and XLM-RoBERTa pre-trained tokenizer to tokenize and pad the training data. As for the test set, I used this test set shared by one of the competitors [2] which translates the entire test set into English. When I evaluate my XLM-RoBERTa model, I can do it directly without translating the test set because it is a cross-lingual model.

Since this is a kernel-only competition, inferencing is done separately on a Kaggle notebook.

4 Models

This section investigates three different models.

4.1 BERT for Text Classification Directly

To adapt pre-trained BERT base model to do text classification, I just need to add a classifier layer on top of the 12 BERT layers. Then, I can fine-tune the hyper-parameters to improve performance.

I will add drop out for BERT base to prevent overfitting. BERT base model produces two outputs: last hidden state and pooled output. After BERT base, I will do a mean pooling and max pooling on the last hidden state from BERT base. Then, I will concatenate them and send it to a linear (output) layer.

This linear (output) layer has size of (768*2,1). 768 is the hidden size of BERT base, which means it has 768 output features, we times it by 2 because we did both max pooling and mean

pooling before. Since it is a binary classification problem, I put 1 to make the output value between 0 and 1.

I will use binary cross entropy with logits loss as our loss function, and Adam with decoupled weight decay[9] as our optimizer.

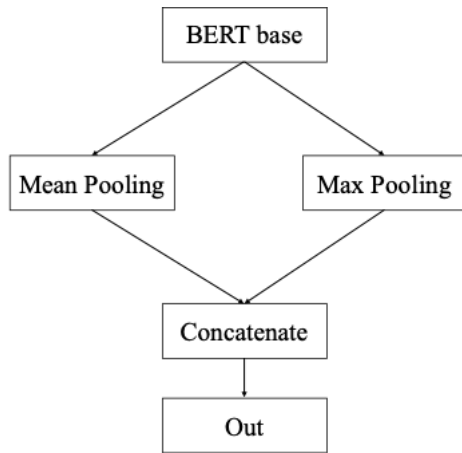


Figure 5: BERT for Text Classification

4.2 BERT as word Embeddings + KimCNN

I will use BERT to tokenize and pad words into word embeddings. Then I will feed those word embeddings to a CNN introduced in this paper [7] by Yoon Kim from New York University in 2014.

Here are the steps of how KimCNN works in details [7]:

- (1) After we get word embeddings from BERT, which are of size [192, 768]. The max sequence length is 192 and the embedding size is 768.
- (2) The model applies convolution on embeddings that uses three different filters of size [2, 768], [3, 768], and [4, 768]. These filters correspond to capturing the bi-gram, tri-gram, and 4-gram relations.
- (3) Apply ReLU to add the ability to model nonlinear problems.
- (4) Apply 1-max pooling to down-sample the output representation of the convolutional layer and to help prevent overfitting.
- (5) Concatenate vectors after 1-max pooling

into a single vector.

- (6) Dropout (to prevent overfitting).

- (7) Apply sigmoid function.

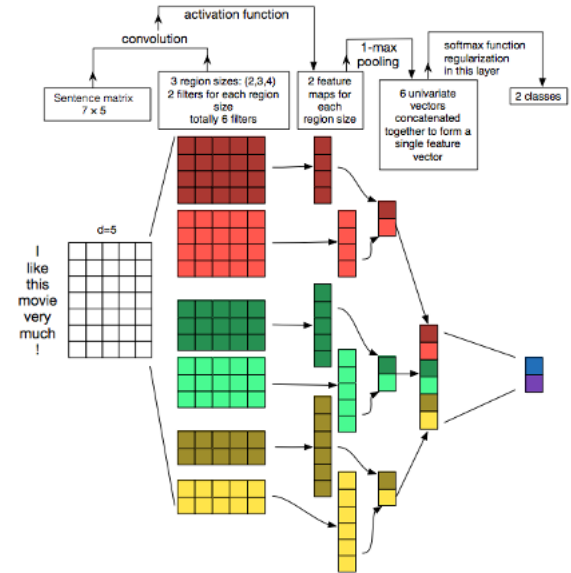


Figure 6: KimCNN [7]

4.3 Cross-lingual model: XLM-RoBERTa

This is very similar to BERT. The pre-trained XLM-RoBERTa has two outputs, last hidden state and pooled output. The pooled output is feed into a linear layer of size (1024, 1). 1024 is the hidden size of XLM-RoBERTa.

5 Experiments

I have fine-tuned BERT and XLM-RoBERTa to compare against CNN. Also, I have trained BERT both on GPU and TPU to compare the performances.

5.1 Fine-tune BERT on GPU

For fine-tuning, I will tune these hyperparameters: learning rate, batch size and number of training epochs.

Since the training dataset is rather large and training on GPU takes a long time, I will tune the hyperparameters on the first 50,000 datapoints.

The BERT paper had a list of hyperparameters that generally work well across all tasks,

which I can use as my starting point to fine-tune [4]: Batch size: 16, 32; Learning rate: $5e-5$, $3e-5$, $2e-5$; Number of epochs: 2, 3, 4.

Learning Rate: Learning rate is perhaps the most critical one of all the hyperparameters. To tune learning rate, I will fix the training batch size to be 64 and the validation batch size to be 8. Training on one epoch. The training and validation losses plot every 100 batches to make the plot less clustered. I will use ROC-AUC score to measure how well our model performs on validation set. Please note that the y-axes of figures below are at different scales.

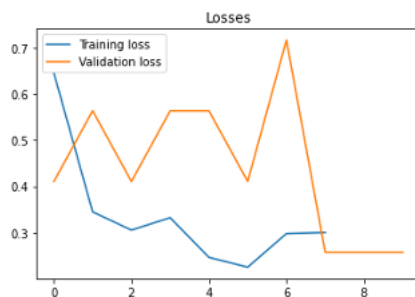


Figure 7: $lr = 5e-3$

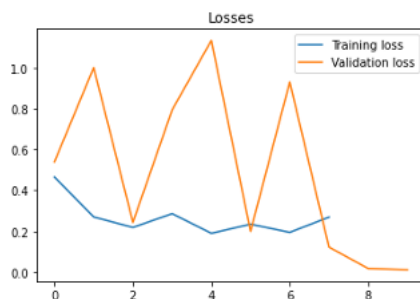


Figure 8: $lr = 5e-5$

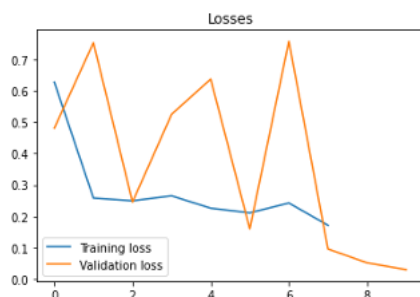


Figure 9: $lr = 2e-5$

As shown by Table 1, with an aggressive learning rate of $5e-3$, the model did poorly. I find that

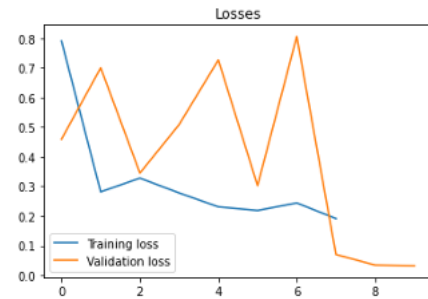


Figure 10: $lr = 5e-6$

Learning Rate	ROC-AUC
$5e-3$	0.53
$5e-5$	0.82
$2e-5$	0.84
$5e-6$	0.79

Table 1: Learning Rate vs. ROC-AUC

a good learning rate for this model should be around $2e-5$, which is suggested by the paper.

Batch Size: I will set learning rate to be $1e-5$, training on 3 epochs to tune batch sizes.

Generally, too large of a batch size will lead to poor generalization, while too small of a batch size will not guarantee convergence to the global optimum. It will bounce around the global optima instead. [10]

I have found that the max training batch size is 64. I tried training batch sizes of 128 and 256, which will blow up GPU memory.

Training	Validation	ROC-AUC
16	2	0.79
32	4	0.81
64	8	0.82

Table 2: Batch Size vs. ROC-AUC

As we can see training batch size of 64 is ideal for this model given my constraints.

Number of Epochs: After trained on 6 epochs with learning rate of $1e-6$, training batch size of 64, validation batch size of 8. We get an ROC-AUC score of 0.83 at the 6th epoch. As the original paper suggested, BERT is far less sensitive to hyperparameter tuning on large data

sets (100k+)[4]. As we can see from the results, once the hyperparameters are set to around the optimal values, it does not affect the model much.

My fine-tuned BERT base model achieved the accuracy of 0.9012 on the test set on Kaggle. This model takes 4.7 hours to train on GPU.

Max Sequence Length	192
Training Batch Size	64
Validation Batch Size	8
Epochs	5
Learning Rate	1e-6

Table 3: Best Performed BERT

5.2 Fine-tune BERT on TPU

I trained my best BERT model on TPU. I used XLA package, which enables PyTorch on Google TPU. To make the training process even faster, I utilized all 8 cores on TPU to parallel processing. I have trained for one epoch to see how fast the speed up on TPU is.

XLA : 1	12m 16s
XLA : 2	12m 12s
XLA : 3	12m 15s
XLA : 4	12m 13s
XLA : 5	12m 18s
XLA : 6	12m 17s
XLA : 7	12m 28s
XLA : 8	12m 19s

Table 4: Training Time on TPU

As shown by Table 4, it only takes 12 minutes to train for one epoch on TPU vs. 56 minutes to train for one epoch on GPU.

5.3 BERT + CNN

Based on my experiments, this model does not perform well on the dataset. The training loss failed to converge to a low rate.

5.4 Fine-tune XLM-RoBERTa

After fine-tuning XLM-RoBERTa, my best performed model achieves 0.9459 on test set on Kaggle.

Max Sequence Length	192
Training Batch Size	16
Validation Batch Size	4
Epochs	2
Learning Rate	5e-6

Table 5: Best Performed XLM-RoBERTa

6 Conclusion

This paper has shown that transfer learning on powerful models like BERT achieves much better result than creating a model from scratch. This is because models like BERT takes a lot of resources to train and it generalize very well on different tasks.

Also, I have found that multi-lingual models performs better than mono-lingual models.

Lastly, TPU is powerful especially comes to training a very large model like BERT.

References

- [1] Jay Alammr. The illustrated transformer.
- [2] Camaro. test-en-df, Mar 2020.
- [3] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. 2019.
- [4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018.
- [5] Jigsaw. Jigsaw llc, Mar 2020.
- [6] Jigsaw. Jigsaw multilingual toxic comment classification, Mar 2020.
- [7] Yoon Kim. Convolutional neural networks for sentence classification. *CoRR*, abs/1408.5882, 2014.
- [8] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692, 2019.
- [9] Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. *CoRR*, abs/1711.05101, 2017.
- [10] PyTorch. Pytorch on xla devices.
- [11] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017.