

Optimizing Inference Performance of Transformers on CPUs

Dave Dice

Alex Kogan

Oracle Labs

Why Transformers?

[cs.CL] 6 Dec 2017

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaizer@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention

Why Transformers?

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova
Google AI Language
{jacobdevlin, mingweichang, kentonl, kristout}@google.com

Abstract

We introduce a new language representation model called BERT, which stands for Bidirectional Encoder Representations from Transformers. Unlike recent language representation models (Peters et al., 2018a; Radford et al., 2018), BERT is designed to pre-train deep bidirectional representations from unlabeled text by jointly conditioning on both

[cs.

There are two existing strategies for applying pre-trained language representations to downstream tasks: *feature-based* and *fine-tuning*. The feature-based approach, such as ELMo (Peters et al., 2018a), uses task-specific architectures that include the pre-trained representations as additional features. The fine-tuning approach, such as the Generative Pre-trained Transformer (OpenAI

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention

You Need

Vi Parmar* Jakob Uszkoreit*
Google Research Google Research
vp@google.com usz@google.com

Lukasz Kaiser*
Google Brain
lukaszkaizer@google.com



Why Transformers?

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova
Google AI Language
{jacobdevlin, mingweichang, kentonl, kristout}@google.

Abstract

We introduce a new language representation model called BERT, which stands for Bidirectional Encoder Representations from Transformers. Unlike recent language representation models (Peters et al., 2018a; Radford et al., 2018), BERT is designed to pre-train deep bidirectional representations from unlabeled text by jointly conditioning on both

[CS:

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

Language model pretraining has led to significant performance gains but careful comparison between different approaches is challenging. Training is computationally expensive, often done on private datasets of different sizes, and, as we will show, hyperparameter choices have significant impact on the final results. We present a replication study of BERT pretraining (Devlin et al., 2019), which includes a careful evaluation of the effects of hyperparameter tuning and training set size. We find that BERT was significantly undertrained and propose an improved recipe for training BERT models, which we call RoBERTa, that can match or exceed the performance of all of the post-BERT methods.

Illia Polosukhin

illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention

Why Transformers?

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova
Google AI Language
{jacobdevlin, mingweichang, kentonl, kristout}@google.

ALBERT: A LITE BERT FOR SELF-SUPERVISED LEARNING OF LANGUAGE REPRESENTATIONS

Zhenzhong Lan¹ Mingda Chen^{2*} Sebastian Goodman¹ Kevin Gimpel²
Piyush Sharma¹ Radu Soricut¹

¹Google Research ²Toyota Technological Institute at Chicago

{lanzhzh, seabass, piyushsharma, rsoricut}@google.com
{mchen, kgimpel}@ttic.edu

ABSTRACT

Increasing model size when pretraining natural language representations often results in improved performance on downstream tasks. However, at some point further model increases become harder due to GPU/TPU memory limitations and longer training times. To address these problems, we present two parameter-efficient, lightweight, domain-independent pretraining methods that train

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

Language model pretraining has led to significant performance gains but careful comparison between different approaches is challenging is computationally expensive on private datasets of different sizes. We will show, hyperparameter tuning will have a significant impact on the final results.

lia Polosukhin
polosukhin@gmail.com

Abstract

Modern neural network architectures for natural language processing are based on complex recurrent or convolutional architectures that often include an encoder and a decoder. The best performing architectures often include an attention mechanism over the input sequence.

Why Transformers?

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova
Google AI Language
{jacobdevlin, mingweichang, kentonl, kristout}@google.

ALBERT: A LITE BERT FOR SELF-SUPERVISED LEARNING OF LANGUAGE REPRESENTATIONS

Zhenzhong Lan¹ Mingda Chen^{2*} Sebastian Goodman¹ Kevin Gimpel²
Piyush Sharma¹ Radu Soricut¹

¹Google Research ²Toyota Technological Institute at Chicago

{lanzhzh, seabass, piyushsharma, rsoricut}@google.com
{mchen, kgimpel}@ttic.edu

ABSTRACT

Increasing model size when pretraining natural language representations often results in improved performance on downstream tasks. However, at some point further model increases become harder due to GPU/TPU memory limitations and longer training times. To address these problems, we present two parameter-efficient algorithms for pretraining natural language representations.

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, danqi, omerlevy, mikelewis}@fb.com

Abstract

Language model pretraining performance between different training is constrained on private context we will show significant improvement in accuracy

Polosukhin

Abstract

ation mode that include the model

Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context

Zihang Dai^{*12}, Zhilin Yang^{*12}, Yiming Yang¹, Jaime Carbonell¹, Quoc V. Le², Ruslan Salakhutdinov¹

¹Carnegie Mellon University, ²Google Brain

{dzihang, zhiliny, yiming, jgc, rsalakhu}@cs.cmu.edu, qvl@google.com

Abstract

Transformers have a potential of learning longer-term dependency, but are limited by a fixed-length context in the setting of language modeling. We propose a novel neural architecture *Transformer-XL* that enables learning

Term Memory (LSTM) networks (Hochreiter and Schmidhuber, 1997), have been a standard solution to language modeling and obtained strong results on multiple benchmarks. Despite the wide adaption, RNNs are difficult to optimize due to gradient vanishing and explosion (Hochreiter and Schmidhuber, 1997).

Why Transformers?

Can be pre-trained on huge amounts of unlabeled data

- E.g., all of Wikipedia, book corpus, etc.

Fine-tuned later to a specific task

- E.g., question-answering, sentiment analysis, text-classification
- with just a small amount of labeled, domain-specific data

Feature millions/billions of parameters

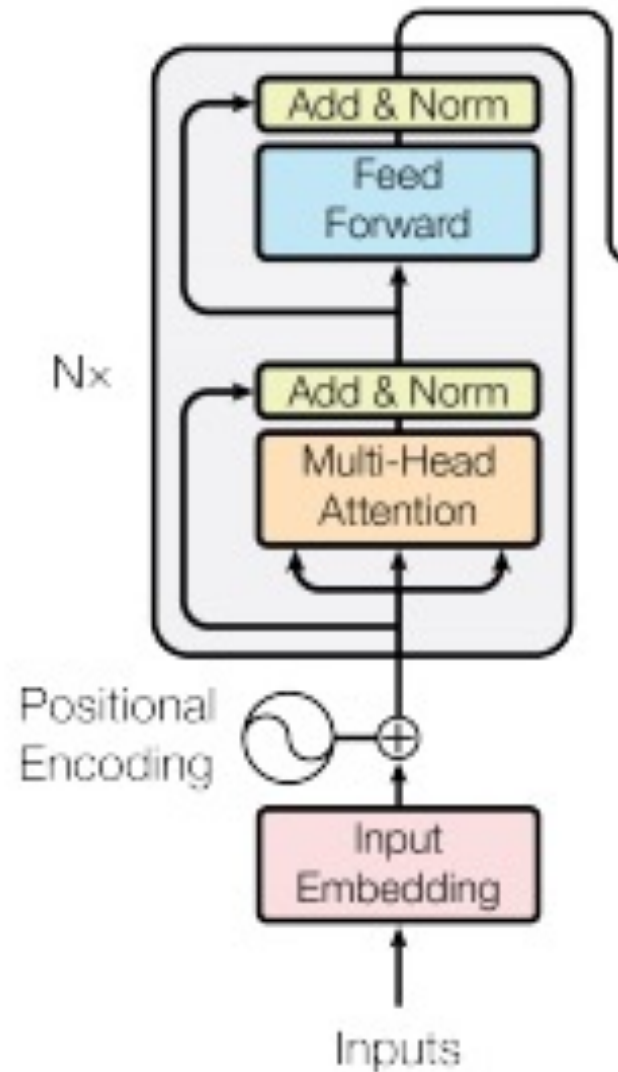
→ Lots of attention on optimizing the training process

- (relatively) less attention on optimizing inference

Why CPUs?

- Commodity hardware
- Cost-effective choice for real-time (aka batch=1) inference
- Preferred choice for inference deployment in industry

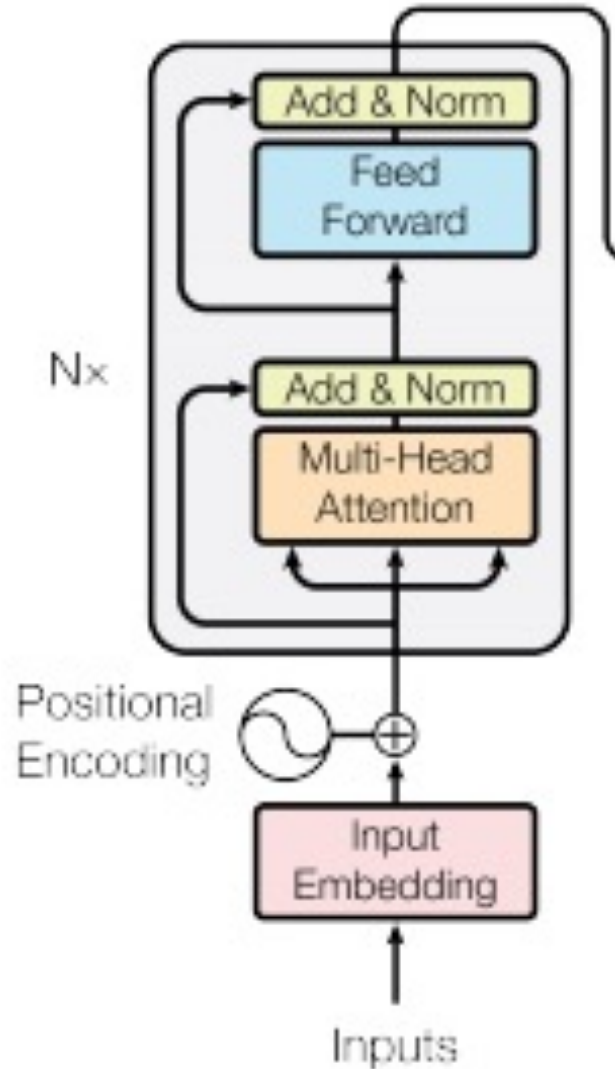
Transformer Architecture



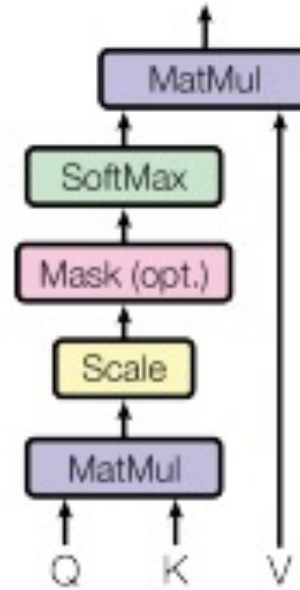
From “Attention is all you need” by Vaswani et al., 2017



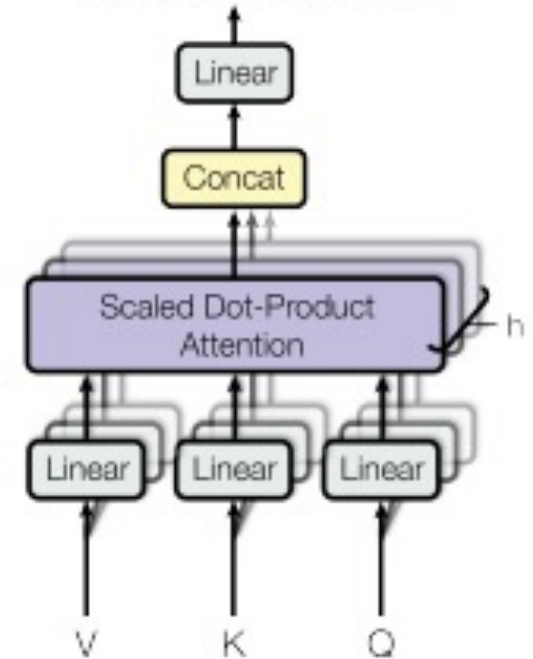
Transformer Architecture



Scaled Dot-Product Attention



Multi-Head Attention



Our setup

OCI BM.Standard2.52 instance

- Intel Skylake-gen processor with 26 hyper-threaded cores
- AVX512 support

Pytorch (v1.6)

oneDNN math library

- enhanced integration of oneDNN with Pytorch

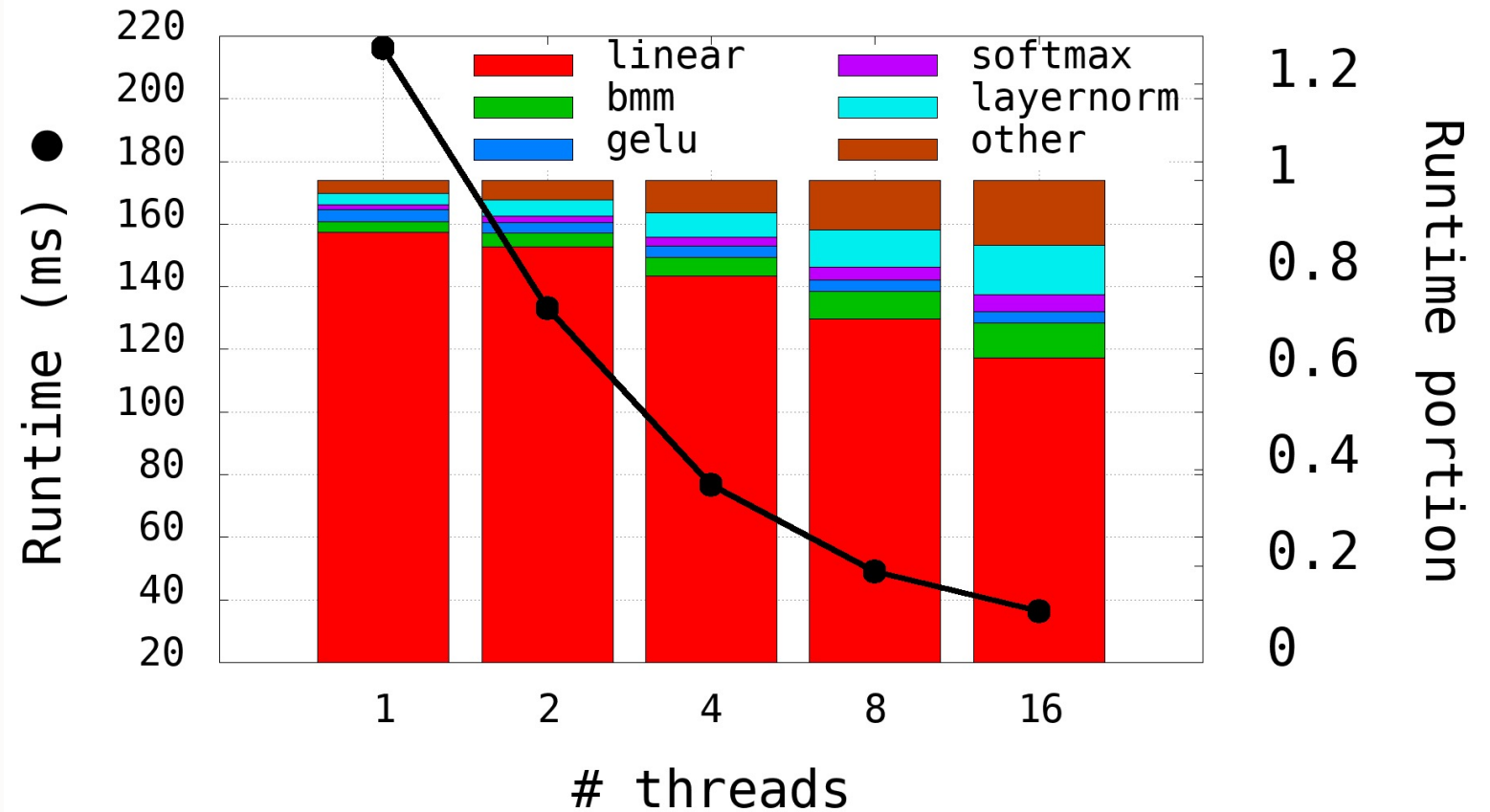
Transformers python package (v3.0.2) from HuggingFace

- state-of-the-art implementation of numerous Transformer-based NLP models
- “BERT-base” model (from Transformers)

Runtime breakdown by modules

- Relatively poor scalability
 - ~5x from 1 to 16 threads
- Matmuls scale better than other ops
 - their share decreases with increase in #threads

Between 66 and 92% of the time is spent in matmul ops



Key observation

The **matmul operation is heavily impacted** not only by the **shape** (dimensions) of the source matrices and available **resources** (#threads), but also by the **form** (transposed/normal) of those matrices

- the form dictates the memory layout (column-major vs. row-major)

A subset of Basic Linear ALgebra (BLAS) functions to perform matrix-matrix multiplication. [More...](#)

Functions

`mkldnn_status_t` MKLDNN_API `mkldnn_sgemm` (char transa, char transb, `mkldnn_dim_t` M, `mkldnn_dim_t` N, `mkldnn_dim_t` K, float alpha, const float *A, `mkldnn_dim_t` lda, const float *B, `mkldnn_dim_t` ldb, float beta, float *C, `mkldnn_dim_t` ldc)
SGEMM performs a matrix-matrix multiplication operation defined as. [More...](#)

Source: https://oneapi-src.github.io/oneDNN/v1.0/group__c__api__blas.html

Key observation

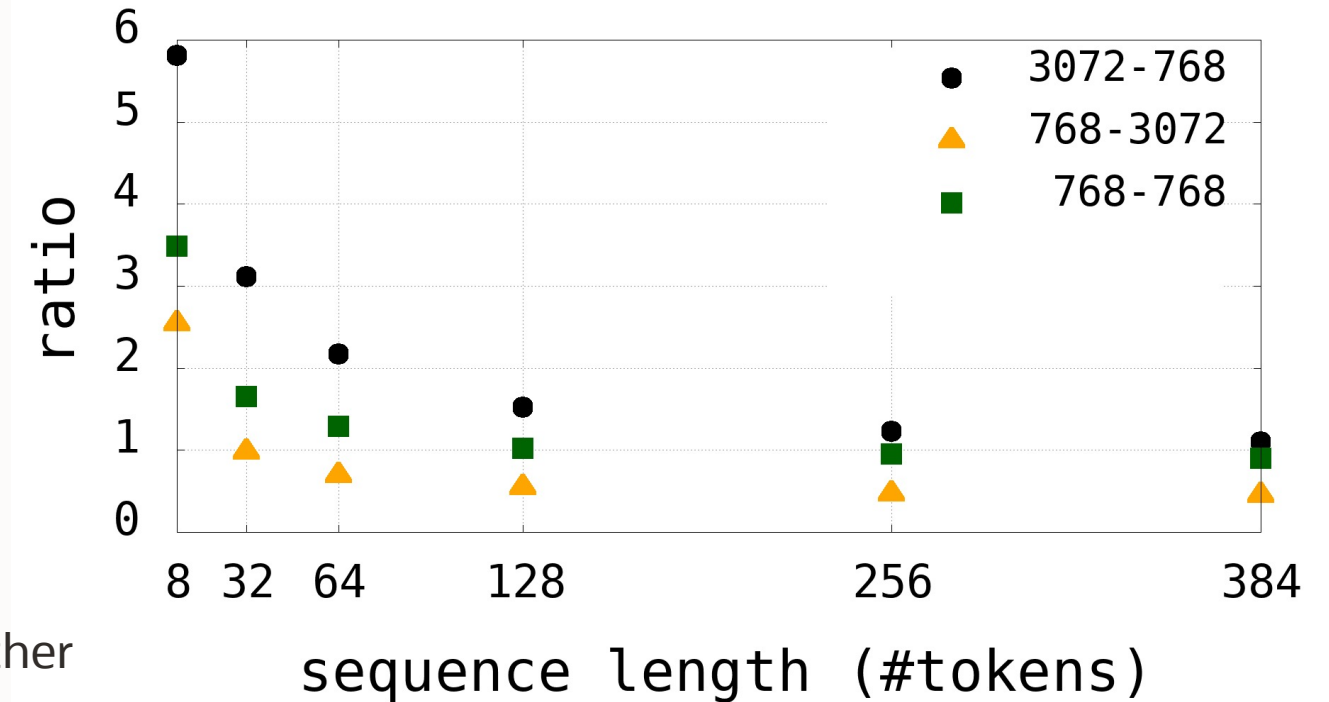
Why is the difference?

- different matmul variants generate different code paths and memory access patterns
- perf. counters suggest faster variants enjoy better L2 cache locality

Why does it matter?

- Pytorch keeps the weights matrix in the Linear module transposed
 - Tensorflow keeps the weights matrix normal
- ➔ Matmuls are always applied in one form or another

☞ But the best approach – adaptivity!



$$\text{ratio} = \frac{A*B, \text{ where } A \text{ and } B \text{ are normal}}{A*B, \text{ where } A \text{ is normal and } B \text{ is transposed}}$$

Adaptive Linear Module Optimization (ALMO)

Augment each Linear module* with `transposeFlags` array of Boolean flags

When creating a Linear module with the shape `[in ,out]`:

- For each `i` in `[0, 9]`:
 - Generate random matrix with the shape `[2i, in]`
 - Measure time to perform matmul with the weight matrix transposed and normal
 - Record the result (0/1) in `transposeFlags[i]`

During inference, with the input of shape `[length, in]`:

- Calculate `s = int(log(length))`
- Based on `transposeFlags[s]`, perform matmul with weights either transposed or not

* `transposeFlags` can be shared among Linear modules of the same shape

Experimental results

| #threads | sequence length | | | | | | | | |
|----------|-----------------|------------------|----------------|----------------|------------------|----------------|----------------|------------------|----------------|
| | 8 | | | 64 | | | 384 | | |
| | onednn base | onednn normal | onednn almo | onednn base | onednn normal | onednn almo | onednn base | onednn normal | onednn almo |
| 1 | 115 | 82 | 79 ± 1 (x1.46) | 216 | 193 | 162 (x1.33) | 884 ± 31 | 1009 | 807 (x1.10) |
| 2 | 83 | 50 | 50 ± 1 (x1.66) | 133 | 105 | 97 ± 1 (x1.37) | 471 ± 29 | 512 | 413 (x1.14) |
| 4 | 51 | 34 | 34 (x1.50) | 76 | 64 | 60 (x1.27) | 259 ± 27 | 285 | 220 (x1.18) |
| 8 | 35 | 27 | 27 (x1.30) | 49 | 45 | 42 (x1.17) | 135 | 148 | 128 (x1.05) |
| 16 | 28 | 24 | 24 (x1.17) | 36 | 34 | 32 (x1.12) | 96 ± 14 | 97 | 82 (x1.17) |

BERT-base inference latency (ms)

Similar results for other models: RoBERTa, DistilBERT



Conclusions

- Faster matmul operations → faster inference
- Adaptive Linear Module Optimization
 - lightweight
 - transparent (no change to the model)
 - does not impact accuracy
 - is not limited to Transformers
- More details in the extended paper (<https://arxiv.org/pdf/2102.06621>):
 - Two more optimizations for the matmul operations
 - Sequential overhead reduction
 - Modified matrix partitioning
 - More experiments, results, etc.