

---

# Efficient Softmax Approximation for Deep Neural Networks with Attention Mechanism

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        There has been a rapid advance of custom hardware (HW) for accelerating the  
2        inference speed of deep neural networks (DNNs). Previously, the softmax layer  
3        was not a main concern of DNN accelerating HW, because its portion is relatively  
4        small in multi-layer perceptron or convolutional neural networks. However, as the  
5        attention mechanisms are widely used in various modern DNNs, a cost-efficient  
6        implementation of softmax layer is becoming very important. In this paper, we  
7        propose two methods to approximate softmax computation, which are based on  
8        the usage of LookUp Tables (LUTs). The required size of LUT is quite small  
9        (about 700 Bytes) because ranges of numerators and denominators of softmax are  
10       stable if normalization is applied to the input. We have validated the proposed  
11       technique over different AI tasks (object detection, machine translation, speech  
12       recognition, semantic equivalence) and DNN models (DETR, Transformer, BERT)  
13       by a variety of benchmarks (COCO17, WMT14, WMT17, GLUE). We showed  
14       that 8-bit approximation allows to obtain acceptable accuracy loss below 1.0%.

## 15    1 Introduction

16    After Vaswani et al. had introduced Transformer model in [27] for machine translation task, the  
17    attention based architecture became popular firstly in Natural Language Processing (NLP) appli-  
18    cations, e.g.: speech recognition [16], [21], [9]; summarization [7]; language understanding [5],  
19    [33], [15], [18]; and video captioning [3]. Recently attention-based models are used in even wider  
20    practical areas including Computer Vision (CV) tasks for object detection [2]; image transformation  
21    [26]; image classification [6]; and even symbolic integration and solving differential equations [14].  
22    Despite the attractiveness of transformer-based models, its direct implementation into the platform  
23    with constrained computational power (e.g., mobile SoC, edge devices) <sup>1</sup> is very challenging due to  
24    big memory footprint and latency.

25    Therefore, model compression techniques such as quantization, and distillation are needed for those  
26    models. Many approaches have been introduced on quantizing matrix multiplication of transformer  
27    architecture. For example, in [22] it was used second order Hessian information, what allows to  
28    significantly compress the size of the model up to  $13\times$  times, while maintaining at most 2.3% of  
29    performance degradation (for the case of ultra-low precision 2-bits quantization). In [35] it was  
30    used a quantization-aware training during the fine-tuning phase of BERT, what allows to compress  
31    BERT model by  $4\times$  (with 8-bit quantization) with minimal accuracy loss (less than 1%). In [1], a  
32    machine language translation model was quantized by 8-bit, while maintaining less than 0.5% drop

---

<sup>1</sup>Recently, interest to the computations performed close to the data sources is growing up aiming to soften the requirements of continuous access to high-speed and high-bandwidth connections. Moreover, often customers wanted to keep their security and privacy, and thus do not want to expose their data to the external clouds [32], [19], [36], [11].

33 in accuracy. Moreover, in [20] it was shown that 8-bit quantized models provide the same or even  
34 higher accuracy as the full-precision models. Most of the above methods consider the quantization of  
35 matrix multiplications operations only. However, as it is shown in [4], [24] in modern DNNs with  
36 attention mechanism (e.g., Transformer, BERT, GPT-x) the softmax function is also used intensively,  
37 especially at the longer sequence lengths, so it is necessary to optimize its for performance.

38 In this paper we propose methods for efficient computation of the softmax layer at the HW accelerator.  
39 The method is based on piece-wise-constant approximation and usage of LUTs. To the best of our  
40 knowledge, it is the first paper where softmax quantization of the models with attention mechanism is  
41 tested and verified on a variety of AI tasks. In Section 2 we show why our research is important and  
42 valuable. In Section 3 we consider the drawbacks of existed softmax approximation methods in the  
43 perspective of HW accelerator, and summarize the differentiation of our methods from the previous  
44 arts. In Section 4 we describe the details of the proposed methods. Section 5 shows the experimental  
45 validation over different models and datasets, and Section 6 concludes the paper.

## 46 2 Background and motivation

47 Modern GPUs are powerful, but big, expensive, and power-hungry. Therefore, alternative HW  
48 accelerators (e.g., NPU) for on-device inference are under active development by different vendors,  
49 especially for Federated Learning and Edge computing. However, such devices mostly are focused on  
50 the acceleration of matrix multiplication operations, and do not include means to compute complex  
51 activation functions. Typically, in such devices the data is sent outside of the accelerator to compute  
52 activations on host CPU. For example, according to the guidelines of Coral (TM), a softmax layer of  
53 DNN model in Edge TPU have to be run on host CPU<sup>2</sup>, what is acceptable for traditional CV tasks  
54 (which are typically uni-directional, have minimum dependencies, and softmax layer is located at the  
55 end of the computational graph of DNN model), however is very inefficient for NLP tasks (which are  
56 typically more complicated with a lot of dependencies and active employment of softmax layer in the  
57 middle of DNN model). In opposite to traditional logic-centric approach, some researches are trying  
58 to perform computation closer to the memory (so called memory-centric approach). For example  
59 in [23], there is shown a DRAM-based AI accelerator. This approach allows significantly speed-up  
60 the overall computation process, but for the computation of the activations the data should also be  
61 moved to host processor, what is an even bigger issue in the DRAM environment.

62 The Eq. (1) from [27] describes how attention is computed in the model. This particular form,  
63 named "scaled dot-product attention" takes the matrix multiplication product of queries and keys  
64 of  $\mathbf{R}^{N \times L \times H}$  as input for the softmax layer where  $N$  means number of heads,  $L$  means sequence  
65 length and  $H$  means hidden size for the case where batch size equals to 1. In other words, performing  
66  $(N \times L \times L)$  softmax operations is required per one attention. Furthermore, encoder in typical  
67 transformer consists of six multi-head attentions which means  $6 \times (N \times L \times L)$  operation is required  
68 for encoder solely. Assuming the number of heads is 8 and sequence length is 128, it already takes  
69 786, 432 operations for softmax of the transformer encoder. This overhead increases as number of  
70 heads and sequence length increases which is typical case for high-performing models.

$$Attention(Q, K, V) = Softmax\left(\frac{QK^T}{\sqrt{d_K}}\right)V \quad (1)$$

71 For example, it requires performing  $12 \times (12 \times 128 \times 128) = 2,359,296$  softmax operations for  
72 one sample inference for typical BERT configuration [5] over sequence of length 128. In the case  
73 when HW accelerator is used for matrix multiplication only, and activation to be computed at the  
74 CPU (what is common case for HW accelerators, optimized for CNN-models), the huge amount of  
75 data must be moved between CPU and the accelerator. Such data movement negatively impacts on  
76 the overall computation time and power consumption, which can be critical for on-device inference.  
77 Therefore, HW accelerator must be able to compute softmax layer without CPU involvement.

## 78 3 Related work and key contributions

79 The common equation to compute softmax function over the input  $x$  is a fraction as shown below:

<sup>2</sup>[https://coral.ai/docs/reference/edgetpu.learn.backprop.softmax\\_regression/](https://coral.ai/docs/reference/edgetpu.learn.backprop.softmax_regression/)

$$\sigma(x_i) = \frac{e^{x_i}}{\sum e^{x_i}} = \frac{e^{x_i - \max(x)}}{\sum e^{x_i - \max(x)}} \quad (2)$$

80 There are different ways to implement it. For example, some approaches straightforwardly compute  
 81 the numerator and denominator firstly, and then a division operation is performed. In such case the  
 82 HW accelerator should contain a divider, what requires additional HW costs and can also cause  
 83 performance degradation, if divider is not fully pipe-lined.

84 In [25] it is proposed to use basic-split calculation method, which allows to split the exponentiation  
 85 calculation of the softmax into several specific basics which are implemented by LUT (ROM). It  
 86 allows to simplify the complexity of hardware and signal propagation delay. However, to recover  
 87 the whole computed value of exponent some additional multiplications are needed. Moreover, to  
 88 obtain the final value of softmax the division is still used. In [30] it is proposed to add threshold  
 89 layers to accelerate the training speed and replace the Euler’s base value with a dynamic base value  
 90 to improve the network accuracy. Such approach allowed to save up to 15% of training model  
 91 convergence time and also increase by 3 to 5% the average accuracy. But during the computation  
 92 of softmax the divider is still used. In [17] the combination of LUT and multi-segment polynomial  
 93 fitting have been used to compute exponential operations of integer and fractional parts in separate.  
 94 In addition, they adopt radix-4 Booth-Wallace based multiplier for computing the whole value  
 95 of exponent, and modified shift-compare divider for computation of the final value of softmax.  
 96 To avoid big area costs for traditional divider, the authors in [8] propose to reduce the operand  
 97 bit-width, and approximate exponential and division operations with cost-effective addition and  
 98 bit shifts operations. In their design they have approximated the division operation in Eq. (2) by  
 99 replacing the denominator with closest  $2^b$  value, where  $b$  is some integer constant. Then division  
 100 is implemented just as simple bit shifts operation. In [24] it is proposed to replace the base as  
 101  $e^x \rightarrow 2^x$ , then all computations are more hardware-friendly, however the division operation is still  
 102 required. Also, to restore accuracy a fine-tuning of the model is needed, what is not applicable  
 103 for post-training quantization paradigm. Although the methods described above are decreasing the  
 104 hardware complexity of softmax computation they all still rely on the division operation.

105 To avoid division operation at all, some other solutions apply the logarithmic transformation to  
 106 the original softmax function, and thus substitute costly division operation by subtraction of the  
 107 logarithm. In [34], for example, it is proposed to use a logarithmic operation implemented as a LUT  
 108 and a subtractor to replace the division operation, what allows to further decrease the complexity of  
 109 hardware, as well critical path of the whole design. In [10] simplified version of Integral Stochastic  
 110 Computation is used in order to build FSM-based exponentiation. Division operation is substituted  
 111 by LUT-based logarithmic operation and subtraction, similarly to [34]. In [31] the authors are further  
 112 developing the method proposed in [34], by applying mathematical transformations and linear fitting.  
 113 After optimization, their final design includes only shift operations, leading one detector, and adders.  
 114 Finally, there are some extreme approximation cases represented in [13] and [28], where logarithmic  
 115 computation and subtraction are skipped at all.

116 Despite its attractiveness, logarithmic transformation approach can be used only in the cases when  
 117 softmax layer is the last layer in DNN and its functionality is simply “scoring” among the candidates  
 118 for classification tasks. However, if softmax layer is used inside of computational graph of DNN (e.g.,  
 119 DNNs with attention-mechanism) then error caused by quantizations will be accumulated drastically,  
 120 directly impacting on the final accuracy. For example, in Table 1 there is shown the averaged accuracy  
 121 drop for DETR models caused by a softmax approximation in uint8 precision by some prior arts.  
 122 As it can be seen from the Table 1, a straightforward usage of Eq.(2) from [31] causes big accuracy  
 123 drop, and even after applying some improvements to the original method (shown as case Eq.(2)+),  
 124 the accuracy drop is still high (2% to 19%). For more details of the prior arts experiments please refer  
 125 to Appendix A.1. However, if for the same conditions we use the method proposed in Section 4.1, we  
 126 can see that accuracy drop reduced by  $\times 4$  to  $\times 20$  times, and it is below 0.5% for plain DETR models  
 127 (no DC5 dilation at the last stage).

128 The work presented in this paper has focused on the development of methods for efficient computation  
 129 of softmax layer during the inference at the edge devices, what usually have limited computational  
 130 power and suffer from constraints of the bandwidth.

131 Previous works for HW accelerator of softmax layer are focused on the logic-centric approach and  
 132 used dedicated hardware for its implementation. In such case the utilization of hardware is low,

Table 1: Averaged accuracy drop by different methods over DETR models (Average Precision), %

METHOD	DETR (R50)	DETR+DC5(R50)	DETR (R101)	DETR+DC5(R101)
EQ.(2) IN [31]	7.20	19.30	10.25	25.37
EQ.(2)+ IN [31]	2.50	12.93	5.38	18.85
SECTION 4.1	<b>0.33</b>	2.92	<b>0.22</b>	2.73

133 performance can be slower, and no reconfigurability is provided. In our paper we have used an  
 134 alternative memory-centric approximate computing approach. It keeps accuracy loss small, while  
 135 allows computing softmax operation with no divider. The size of the required memory (i.e., LUT) is  
 136 reasonably small and can be reconfigured on demand.

137 The methods proposed in the paper contribute to building the alternative concept of hardware  
 138 architecture to accelerate essential operations for AI applications, especially for on-device inference.  
 139 To summarize, we have three-fold difference from the previous works:

- 140 • Applicability of our methods to DNN with **attention mechanism** is experimentally proven  
 141 over variety of the models for different AI applications. All previous methods were used  
 142 only for the cases when softmax is the last layer in DNN, and is used for “scoring”.
- 143 • **No divider** is needed to fully implement the method. Moreover, for 2D LUT method even  
 144 multiplier is not needed. Thus, hardware overhead is minimal, and is almost free if used in  
 145 the DRAM-based AI accelerator.
- 146 • Our solutions utilize **integer precision**, what makes it compatible with traditional HW  
 147 accelerators used for matrix multiplication, and simplify the integration of methods into full  
 148 system (all prior methods are based on a fixed point precision).

## 149 4 Proposed methods

150 In this paper we use memory-centric approach to build the accelerator for softmax computation  
 151 in hardware platform with limited resources. We propose two LUT-based methods for efficient  
 152 computation, which provide high performance and do not require a divider. The details of the  
 153 methods are described below and appropriate software models are shown in Appendix A.2.

### 154 4.1 Normalization of reciprocal exponentiation

155 In this subsection we consider the method, which is based on the normalization of reciprocal  
 156 exponentiation, and hereafter we call it REXP for short.

157 The original reciprocal exponentiation method was proposed in [28], where they used the inverse  
 158 way of max-normalization and the reciprocal of exponential function as below:

$$\sigma^*(x_i) = \frac{1}{e^{\max(x) - x_i}} \quad (3)$$

159 And thus, the final value of softmax can be obtained by reading from a simple LUT-table. Content of  
 160 LUT is computed as shown below:

$$LUT_{1/e}[i] = \left\lceil \frac{1}{e^i} \cdot (2^w - 1) \right\rceil, \forall i = 0, 1, \dots, x_q + 1 \quad (4)$$

161 where  $w$  is a number of bits for quantization, and  $x_q = \lceil \ln(2^w - 1) \rceil$  is an efficient quantization  
 162 boundary.

163 In addition to very low computational complexity, this method has other desired properties [28]:

- 164 • it is positive ( $\frac{1}{e^x} > 0 \forall x \in (-\infty, +\infty)$ );

- 165 • bounded and stable ( $\frac{1}{e^{\max(x)-x_i}} \in (0, 1]$ );
- 166 • and nonlinear ( $\frac{1}{e^{\alpha x}} \neq \alpha \frac{1}{e^x}$ ).

167 But due to its aggressive approximation nature, it can be applied only to simple CV tasks, and if  
 168 used for attention-based DNN models causes the explosion of accuracy drop (see Appendix A.1 for  
 169 details). Thus, in this paper we further develop that method to be applicable for wider class of DNN  
 170 models.

171 During our initial investigations, we have noticed that method described in Eq.( 3) is just scaled  
 172 version of real softmax. So, we proposed to normalize it with some probability density function  
 173 (PDF) scale, such that  $\int PDF = 1$ . However, if used straight-forwardly, it would need to involve a  
 174 division operation, what is strongly un-desirable for devices with constrained computational power.  
 175 Therefore, instead of dividing, we propose to substitute division by multiplication with some PDF  
 176 normalizing constant as below:

$$\sigma(x_i) = \frac{\sigma^*(x_i)}{PDF_{norm}} \rightarrow \sigma^*(x_i) \cdot \alpha \quad (5)$$

177 where  $\alpha = e^{-\ln(\Sigma\sigma^*(x_i))}$  is PDF normalizing constant.

178 Then final equation to compute softmax approximation by proposed REXP method is shown below:

$$\sigma(x_i) = \frac{e^{-\ln(\Sigma\sigma^*(x_i))}}{e^{\max(x)-x_i}} = \frac{1}{e^{\max(x)-x_i}} \cdot e^{-\ln(\Sigma\sigma^*(x_i))} \quad (6)$$

179 Thus, to compute the softmax value it requires just two LUTs of considerably small size, where  
 180 content of the first LUT is computed accordingly to Eq.(4), and the second LUT values can be  
 181 computed as below:

$$LUT_\alpha[j] = \left[ \frac{1}{j} \cdot (2^w - 1) \right], \forall j = 0, 1, \dots, x_s - 1 \quad (7)$$

182 where  $j = \Sigma\sigma^*(x_i)$ ,  $x_s$  is selected quantization boundary, and  $LUT_\alpha[x_s] = 0$ .

## 183 4.2 2-Dimensional LUT

184 In this subsection we propose another method which is based on the substitution of a division  
 185 operation in Eq.( 2) by 2-Dimensional (2D) LUT to speed-up and simplify the computation, while  
 186 maintaining accuracy even for attention-based DNN models. Hereafter we will refer to this method  
 187 as 2D LUT.

188 For this purpose, we have started with the estimation of distributions of  $e^x$  and  $\Sigma e^x$  terms for typical  
 189 inference runs. Our investigation showed that if max-based normalization is applied to the input  
 190 values (i.e.,  $x \rightarrow (x - \max(x))$ ), the distribution of  $e^x$  is stable within range  $e^x \in (0, 1]$  regardless of  
 191 the input values, and range of  $\Sigma e^x$  term depends on the length of the input  $x$ . Thus, it allows us to  
 192 have stable computation even within small size of LUT.

193 Generic architecture and concept of the proposed method for efficient softmax implementation as  
 194 2D LUT is shown in Figure 1. There are two LUTs used: 1D LUT for approximation of  $e^x$  values,  
 195 and 2D LUT for storing softmax output values dependent on the values of numerator  $e^x$  (used as  
 196 the 1-st index in the table), and denominator  $\Sigma e^x$  (used as the 2-nd index) of Eq.(2). As it can be  
 197 seen from Figure 1(right), to calculate the indexes for corresponded value in 2D LUT table, only  
 198 most-significant bits (MSB) are needed. Thus, the simplest hardware realization can be done within  
 199 wiring only (when MSB bits are directly connected to the appropriate address selectors)<sup>3</sup>. Also, the  
 200 proposed method can be easily modified to the case where, 1-st index of 2D LUT table is calculated  
 201 not from  $e^x$  but directly from input  $x$ . In such case there is no need to store intermediate values of  $e^x$ .

202 While the content of 1D LUT for approximation of  $e^x$  values is straightforward, 2D LUT contains  
 203 the family of linear approximations where each row contains the softmax output scaled according to

<sup>3</sup>Other hardware realization are also possible, but not considered here for simplicity of the explanation.

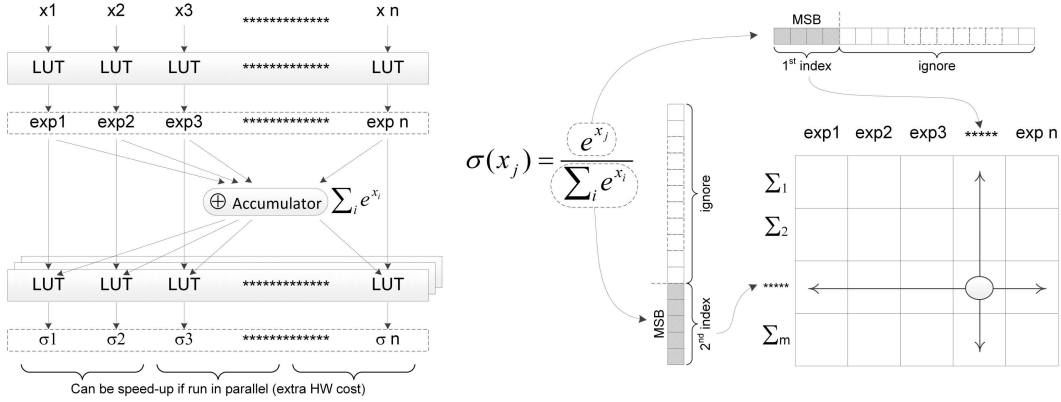


Figure 1: Generic concept of the proposed 2D LUT method (left). Reading softmax output value from pre-computed 2D LUT(right). Computational flow consists from two steps: a) obtaining of  $e^{x_i}$  values by reading from 1D LUT and accumulation of  $\Sigma e^x$  term, b) obtaining  $\sigma(x_i)$  values by reading from 2D LUT.

204  $\Sigma e^x$  term as shown in Eq.(8). The indexes of LUT are computed according to Eq.(9) and Eq.(10),  $w$   
 205 means the number of bits for the value in selected precision.

$$LUT_{\sigma}[i][j] = \left\lfloor \frac{i \cdot scale_{e^x}}{j \cdot scale_{\Sigma}} \cdot (2^w - 1) \right\rfloor \quad (8)$$

206 where

$$i = 0, \dots, \left\lfloor \frac{\max(e^x)}{scale_{e^x}} \right\rfloor \quad (9)$$

$$j = 1, \dots, \left\lfloor \frac{\max(\Sigma e^x)}{scale_{\Sigma}} \right\rfloor \quad (10)$$

207 Since  $x \rightarrow (x - \max(x))$  normalization was used, so  $\max(e^x) = 1.0$ . Therefore,  $scale_{e^x}$  factor  
 208 allows to define the number of columns in LUT to make it small enough for practical applications. In  
 209 our experiment we have selected  $scale_{e^x} = 0.1$  for all precisions, what allows us to reduce the size  
 210 of LUT significantly (i.e.,  $i = 0, \dots, 10$  for all versions of  $LUT_{\sigma}$ ). The value of  $\max(\Sigma e^x)$  depends  
 211 on the distribution of input values. Our experiments showed that  $\max(\Sigma e^x) = 60$  is big enough for  
 212 the tested NLP applications. We also selected  $scale_{\Sigma} = 1.0$  for simplicity of the computations. Thus,  
 213 finally, those parameters give us  $LUT_{\sigma}$  of typical size  $11 \times 60$ .

## 214 5 Experimental validation

215 To validate the proposed methods and check how well they generalize we have conducted several  
 216 experiments with different models (DETR, Transformer, and BERT) for different applications (object  
 217 detection, machine translation, sentiment analysis, and semantic equivalence) over variety of datasets.  
 218 In all those experiments we have used available pre-trained models, where we applied dynamic post-  
 219 training quantization (hereafter we referred to quantized models as PTQ-D)<sup>4</sup>. Then we substituted a  
 220 conventional softmax layer in quantized models with the LUT-based computation as described in  
 221 Section 4. We did not consider any retraining or fine-tuning of the models after quantization, and  
 222 the same off-line generated LUTs were used among all models. Our code allows to select LUTs  
 223 with different precision from int16 down to uint2, what allows to analyze the sensitivity of the model  
 224 to softmax approximation even for ultra-low 2-bits quantization. The details of experiments are  
 225 described below, and results are summarized in Figure 2, Figure 3, and Table 2. For more details

<sup>4</sup>See more details in Appendix A.3.

226 please refer to Table 6, and Table 7 in Appendix. As it can be seen from figures, proposed LUT-based  
 227 softmax computation methods maintain accuracy drop below 1.0% down to 8-bit quantization for all  
 228 NLP and DETR (no DC5) models.

## 229 5.1 Object detection

230 For our first experiments we have used DEtection TRansofmer (DETR) models for object detection [2],  
 231 with available pre-trained models<sup>5</sup>. As it can be seen from Table 6, we were able to reproduce the  
 232 same results for original FP32 reference model over COCO dataset. We have used the same IoU  
 233 metric by Average Precision (AP) as in Table 1 in [2]. Then we run a bunch of experiments to  
 234 check how accuracy of object detection will be decreased due to PTQ-D quantization and LUT-based  
 235 approximation as proposed in REXP method (see Section 4.1). Table 5 in Appendix shows the LUTs  
 236 size for several pre-selected cases in int16 and uint8 precision. There are three cases selected which  
 237 are different in the size of  $LUT_{\alpha}$ : it is  $1 \times 256$  for case 1,  $1 \times 320$  for case 2, and  $1 \times 512$  for case 3.

238 Analysis of Figure 2 shows that accuracy drop caused by application of softmax approximation is  
 239 small ( $< 1\%$ ) and acceptable for plain DETR models (no DC5 used). Bigger accuracy drop for +DC5  
 240 cases is caused by the bigger size of self-attentions of the encoder (see details in Section 5.3). We  
 241 expect that increasing size of LUTs will help to solve this issue. The behavior of average recall values  
 242 is similar to average precision values.

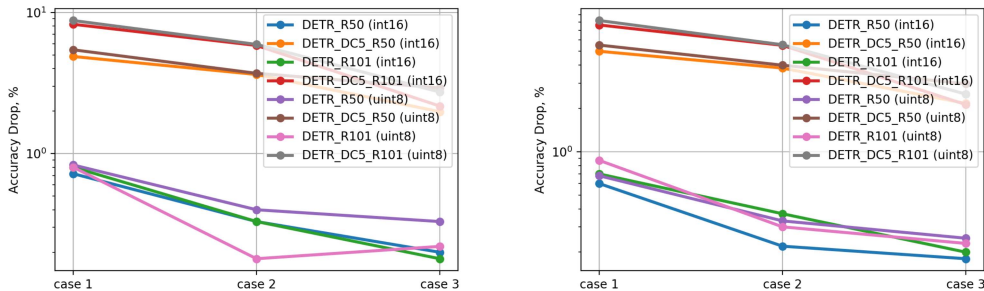


Figure 2: DETR averaged accuracy drop of PTQ-D models with softmax approximations vs. original FP32 models: average precision (left) and average recall (right). As it can be seen from the figure, for DETR models without dilation at the last stage (no DC5) the accuracy drop for all cases is below 1% and shows very similar behavior.

## 243 5.2 NLP tasks

244 Next, we have validated the proposed methods by experimenting with several NLP tasks. Similarly,  
 245 to DETR case, Table 8 in Appendix shows the LUTs size for several pre-selected cases for those  
 246 experiments. In Table 2 below there are accumulated values from the experiments with NLP models.  
 247 The bold values in the table shows highest values per model per method after applying quantization  
 248 and softmax approximation. As it follows from the analysis of experiment results, about 700 Bytes for  
 249 2D LUT method, and up to 50 Bytes for REXP method would be enough for practical applications.

### 250 5.2.1 Machine Translation

251 Among NLP tasks we have started with machine translation. For our experiments we have used  
 252 transformer-base model for En-Ge translation [12] from OpenNMT library, with available pre-trained  
 253 model<sup>6</sup>, configured to replicate the results from original paper. To avoid dependency of the evaluation  
 254 results on the selected tokenization scheme, we have used `spm_decode`<sup>7</sup> to detokenize the output of  
 255 translation, and then applied `multi-bleu.perl` script<sup>8</sup> to calculate BLEU score.

<sup>5</sup><https://github.com/facebookresearch/detr>

<sup>6</sup><https://opennmt.net/Models-py/>

<sup>7</sup><https://github.com/google/sentencepiece>

<sup>8</sup><https://github.com/moses-smt/mosesdecoder>

Table 2: Experimental validation over different NLP models and datasets

PRECISION	TRANSFORMER				BERT			
	2D LUT		REXP		2D LUT		REXP	
	WMT 2014	WMT 2017	WMT 2014	WMT 2017	SST-2	MRPC	SST-2	MRPC
	(BLEU)	(BLEU)	(BLEU)	(BLEU)	(%)	(F1)	(%)	(F1)
FP32	26.98	28.09	26.98	28.09	92.32	90.19	92.32	90.19
PTQ-D	26.86	27.95	26.86	27.95	91.74	89.53	91.74	89.53
INT16	<b>26.87</b>	<b>28.02</b>	<b>26.89</b>	27.64	<b>91.63</b>	<b>89.50</b>	<b>91.74</b>	89.26
UINT8	26.76	27.9	26.8	<b>27.66</b>	<b>91.63</b>	89.35	91.17	<b>89.34</b>
UINT4	26.26	27.43	26.68	28.02	91.40	88.01	91.17	88.77
UINT2	24.42	25.06	25.29	25.86	89.22	56.67	91.63	86.12

256 Thus, as it can be seen from Table 2, we were able to reproduce the same BLEU score for FP32  
 257 reference model as in original model. Then we run several experiments to check how accuracy of  
 258 the translation will be changed due to LUT-based quantization in different precisions, and we can  
 259 confirm that down up to 8-bit quantization the deviation of BLEU score from reference is small for  
 260 both datasets ( $< 0.5\%$ ). Also, if we consider impact of the proposed methods only, then we can  
 261 see that accuracy drop is much smaller, and sometimes even recovers vs. PTQ-D quantization (see  
 262 Figure 3 (right)).

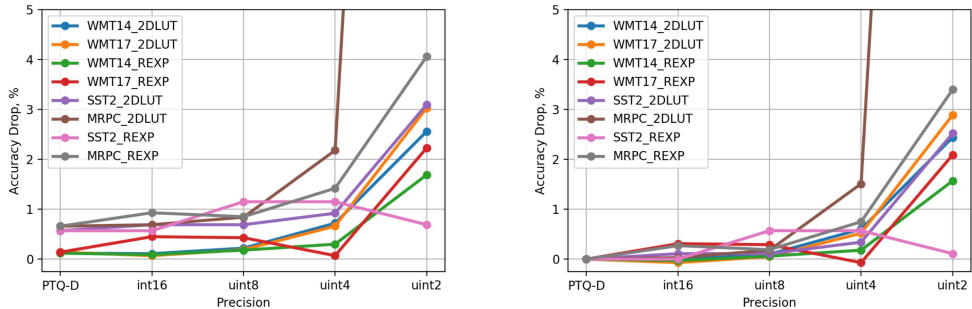


Figure 3: Accuracy drop for NLP experiments: PTQ-D models + softmax approximations vs. FP32 models (left), and PTQ-D models + softmax approximations vs. plain PTQ-D models (right). As it can be seen from the figure, down to uint8 precision the accuracy drop for all cases is below 1% and shows very similar behavior. This confirm very good generalization of the proposed method over different models and applications.

263 **5.2.2 Sentiment analysis**

264 To extend the variety of NLP applications, we also tested the same LUTs with BERT model [5].  
 265 We have used sentiment analysis task from GLUE benchmark [29] to test the model. We have used  
 266 huggingface library <sup>9</sup>, and trained the model with the hyper-parameters described in <sup>10</sup>. The results  
 267 of our experiments showed, that similarly to machine translation, the impact of proposed method  
 268 (softmax layer approximation by LUTs) is smaller vs. accuracy drop caused by PTQ-D quantization  
 269 (see Figure 3).

270 **5.2.3 Semantic equivalence**

271 For semantic equivalence test we used The Microsoft Research Paraphrase Corpus (MRPC) <sup>11</sup> in  
 272 GLUE benchmark. As the classes are imbalanced (68% positive, 32% negative), we follow the

<sup>9</sup><https://github.com/huggingface/transformers>  
<sup>10</sup><https://github.com/google-research/bert>  
<sup>11</sup><https://www.microsoft.com/en-us/download/details.aspx?id=52398>



273 common practice and used F1 score as a metric. We have used `huggingface` library and followed  
 274 the guidelines from PyTorch tutorial <sup>12</sup> to obtain PTQ-D quantized model. Then, similarly to previous  
 275 tests we have substituted a conventional softmax layer with the proposed LUT-based methods. The  
 276 results of our experiments showed the similar trend with sentiment analysis test.

### 277 5.3 Ablation study of DETR models experiment

278 As it is stated in [2], to increase the feature resolution for small objects, a dilation to the last stage of  
 279 the backbone was added (+DC5 cases of DETR models). This modification increases the cost in the  
 280 self-attentions of the encoder, leading to an overall  $\times 2$  increase in computational cost. Such changes  
 281 also reflect on the properties of softmax factors. In Figure 4 there are shown the histogram of  $\Sigma e^x$   
 282 values distributions for the first 200 tensors of DETR model run for `bins = 50`, `range = (0, 500)`.  
 283 As it can be seen from the figure, the distribution of DETR+DC5 (R50) variant is more right-tailed,  
 284 due to the bigger number of high-magnitude values. This causes the bigger accuracy drop when  
 285 LUT-based quantization method is used, due to the lack of the discrepancy for those values. Thus,  
 286 for such models (DETR with added dilation at the last stage) the accuracy of object detection after  
 287 application of the proposed method can be limited. However, as we can see from Figure 2) increasing  
 288 of the size of  $LUT_\alpha$  from 256 Bytes to 512 Bytes allows to decrease the accuracy drop from 9% to  
 289 3% for DETR+DC5 (R101) unit8 case. Thus, we expect that further increasing the size of LUTs will  
 290 help to obtain even more accurate results for DETR models with dilation.

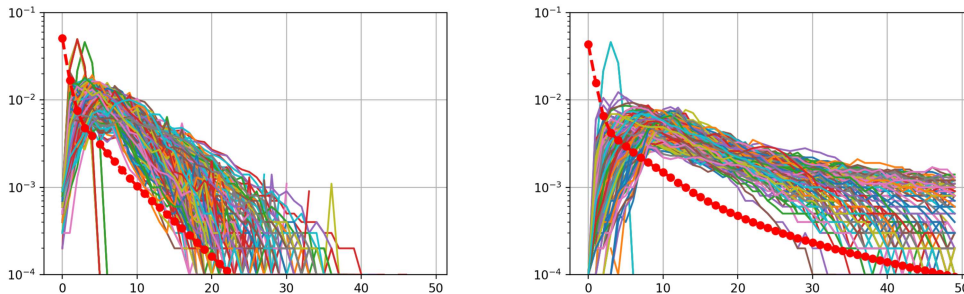


Figure 4: Histogram of  $\Sigma e^x$  values distributions for DETR model variants: plain DETR (R50) (left), and with dilation DETR+DC5 (R50) (right). Red dot line represents the average of the all values per one run of the inference of DETR model. It is clearly seen from the figure that distribution of DETR+DC5 (R50) variant is more flat, having more high-magnitude values. This causes the bigger accuracy drop for the quantized model due to lack of the discrepancy for those values.

## 291 6 Conclusion

292 In this paper two alternative methods for efficient softmax computation for DNN models with  
 293 attention mechanism are proposed. The methods are memory-centric in contrast to known logic-  
 294 centric approach and are based on the usage of LUTs for reading of the pre-computed values, instead  
 295 of the direct computation. Thus, it allows to build the HW accelerator without usage of costly and  
 296 power-hungry divider. In turn, it allows to decrease the power consumption and latency of the whole  
 297 inference, what is crucial for edge computing. All results obtained in the paper were validated  
 298 over different AI tasks (object detection, machine translation, sentiment analysis, and semantic  
 299 equivalence) and models (DETR, Transformer, BERT) by variety of benchmarks (COCO2017,  
 300 WMT14, WMT17, GLUE), showing acceptable accuracy and good generalization of the proposed  
 301 methods.

<sup>12</sup>[https://pytorch.org/tutorials/intermediate/dynamic\\_quantization\\_bert\\_tutorial.html](https://pytorch.org/tutorials/intermediate/dynamic_quantization_bert_tutorial.html)

302 **References**

- 303 [1] Aishwarya Bhandare, Vamsi Sripathi, Deepthi Karkada, Vivek Menon, Sun Choi, Kushal Datta,  
304 and Vikram Saletore. Efficient 8-bit quantization of transformer neural machine language  
305 translation model. *CoRR*, abs/1906.00532, 2019.
- 306 [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and  
307 Sergey Zagoruyko. End-to-end object detection with transformers, 2020.
- 308 [3] Ming Chen, Yingming Li, Zhongfei Zhang, and Siyu Huang. Tvt: Two-view transformer  
309 network for video captioning. In Jun Zhu and Ichiro Takeuchi, editors, *Proceedings of The  
310 10th Asian Conference on Machine Learning*, volume 95 of *Proceedings of Machine Learning  
311 Research*, pages 847–862. PMLR, 14–16 Nov 2018.
- 312 [4] Jacek Czaja, Michal Gallus, Tomasz Patejko, and Jian Tang. Softmax optimizations for intel  
313 xeon processor-based platforms. *CoRR*, abs/1904.12380, 2019.
- 314 [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of  
315 deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018.
- 316 [6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
317 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,  
318 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image  
319 recognition at scale, 2020.
- 320 [7] Sebastian Gehrmann, Yuntian Deng, and Alexander Rush. Bottom-up abstractive summarization.  
321 In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*,  
322 pages 4098–4109, 2018.
- 323 [8] Xue Geng, Jie Lin, Bin Zhao, Anmin Kong, Mohamed M. Sabry Aly, and Vijay Chandrasekhar.  
324 Hardware-aware softmax approximation for deep neural networks. In C.V. Jawahar, Hongdong  
325 Li, Greg Mori, and Konrad Schindler, editors, *Computer Vision – ACCV 2018*, pages 107–122,  
326 Cham, 2019. Springer International Publishing.
- 327 [9] Yanzhang He, Tara N. Sainath, Rohit Prabhavalkar, Ian McGraw, Raziell Alvarez, Ding Zhao,  
328 David Rybach, Anjuli Kannan, Yonghui Wu, Ruoming Pang, Qiao Liang, Deepti Bhatia, Yuan  
329 Shangguan, Bo Li, Golan Pundak, Khe Chai Sim, Tom Bagby, Shuo yiin Chang, Kanishka Rao,  
330 and Alexander Gruenstein. Streaming end-to-end speech recognition for mobile devices, 2018.
- 331 [10] R. Hu, B. Tian, S. Yin, and S. Wei. Efficient hardware architecture of softmax layer in deep  
332 neural network. In *2018 IEEE 23rd International Conference on Digital Signal Processing  
333 (DSP)*, pages 1–5, Nov 2018.
- 334 [11] Andrey Ignatov, Radu Timofte, William Chou, Ke Wang, Max Wu, Tim Hartley, and Luc Van  
335 Gool. AI benchmark: Running deep neural networks on android smartphones. *CoRR*,  
336 abs/1810.01109, 2018.
- 337 [12] Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander M. Rush. OpenNMT:  
338 Open-source toolkit for neural machine translation. In *Proc. ACL*, 2017.
- 339 [13] Ioannis Kouretas and Vassilis Paliouras. Hardware implementation of a softmax-like function  
340 for deep learning. *Technologies*, 8(3), 2020.
- 341 [14] Guillaume Lample and François Charton. Deep learning for symbolic mathematics. *CoRR*,  
342 abs/1912.01412, 2019.
- 343 [15] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu  
344 Soricut. Albert: A lite bert for self-supervised learning of language representations, 2019.
- 345 [16] Bo Li, Shuo yiin Chang, Tara N. Sainath, Ruoming Pang, Yanzhang He, Trevor Strohman, and  
346 Yonghui Wu. Towards fast and accurate streaming end-to-end asr, 2020.
- 347 [17] Z. Li, H. Li, X. Jiang, B. Chen, Y. Zhang, and G. Du. Efficient fpga implementation of  
348 softmax function for dnn applications. In *2018 12th IEEE International Conference on Anti-  
349 counterfeiting, Security, and Identification (ASID)*, pages 212–216, Nov 2018.

- 350 [18] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy,  
351 Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT  
352 pretraining approach. *CoRR*, abs/1907.11692, 2019.
- 353 [19] M. G. Sarwar Murshed, Christopher Murphy, Daqing Hou, Nazar Khan, Ganesh Anantha-  
354 narayanan, and Faraz Hussain. Machine learning at the network edge: A survey, 2019.
- 355 [20] Gabriele Prato, Ella Charlaix, and Mehdi Rezagholizadeh. Fully quantized transformer for  
356 improved translation, 2019.
- 357 [21] Tara N. Sainath, Yanzhang He, Bo Li, Arun Narayanan, Ruoming Pang, Antoine Bruguier, Shuo  
358 yiin Chang, Wei Li, Raziq Alvarez, Zhifeng Chen, Chung-Cheng Chiu, David Garcia, Alex  
359 Gruenstein, Ke Hu, Minh Jin, Anjali Kannan, Qiao Liang, Ian McGraw, Cal Peyser, Rohit  
360 Prabhavalkar, Golan Pundak, David Rybach, Yuan Shangguan, Yash Sheth, Trevor Strohman,  
361 Mirko Visontai, Yonghui Wu, Yu Zhang, and Ding Zhao. A streaming on-device end-to-end  
362 model surpassing server-side conventional model quality and latency, 2020.
- 363 [22] Sheng Shen, Zhen Dong, Jiayu Ye, Linjian Ma, Zhewei Yao, Amir Gholami, Michael W.  
364 Mahoney, and Kurt Keutzer. Q-bert: Hessian based ultra low precision quantization of bert,  
365 2019.
- 366 [23] Hyunsung Shin, Dongyoung Kim, Eunhyeok Park, Sungho Park, Yongsik Park, and Sungjoo  
367 Yoo. Mcdram: Low latency and energy-efficient matrix computations in dram. *IEEE Trans-*  
368 *actions on Computer-Aided Design of Integrated Circuits and Systems*, 37(11):2613–2622,  
369 2018.
- 370 [24] Jacob R. Stevens, Rangharajan Venkatesan, Steve Dai, Bruce Khailany, and Anand Raghu-  
371 nathan. Softmax: Hardware/software co-design of an efficient softmax for transformers.  
372 *CoRR*, abs/2103.09301, 2021.
- 373 [25] Q. Sun, Z. Di, Z. Lv, F. Song, Q. Xiang, Q. Feng, Y. Fan, X. Yu, and W. Wang. A high speed  
374 softmax vlsi architecture based on basic-split. In *2018 14th IEEE International Conference on*  
375 *Solid-State and Integrated Circuit Technology (ICSICT)*, pages 1–3, Oct 2018.
- 376 [26] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and  
377 Hervé Jégou. Training data-efficient image transformers distillation through attention, 2021.
- 378 [27] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
379 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017.
- 380 [28] Ihor Vasylytsov and Wooseok Chang. *Lightweight Approximation of Softmax Layer for On-Device*  
381 *Inference*. Springer International Publishing, 2021.
- 382 [29] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman.  
383 GLUE: A multi-task benchmark and analysis platform for natural language understanding.  
384 *CoRR*, abs/1804.07461, 2018.
- 385 [30] K. Wang, Y. Huang, Y. Ho, and W. Fang. A customized convolutional neural network design  
386 using improved softmax layer for real-time human emotion recognition. In *2019 IEEE Inter-*  
387 *national Conference on Artificial Intelligence Circuits and Systems (AICAS)*, pages 102–106,  
388 March 2019.
- 389 [31] M. Wang, S. Lu, D. Zhu, J. Lin, and Z. Wang. A high-speed and low-complexity architecture  
390 for softmax function in deep learning. In *2018 IEEE Asia Pacific Conference on Circuits and*  
391 *Systems (APCCAS)*, pages 223–226, Oct 2018.
- 392 [32] Siqi Wang, Anuj Pathania, and Tulika Mitra. Neural network inference on mobile socs. *IEEE*  
393 *Design Test*, page 1–1, 2020.
- 394 [33] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell, Ruslan Salakhutdinov, and  
395 Quoc V. Le. Xlnet: Generalized autoregressive pretraining for language understanding. *CoRR*,  
396 abs/1906.08237, 2019.

- 397 [34] B. Yuan. Efficient hardware architecture of softmax layer in deep neural network. In *2016 29th*  
 398 *IEEE International System-on-Chip Conference (SOCC)*, pages 323–326, Sep. 2016.
- 399 [35] Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. Q8bert: Quantized 8bit bert,  
 400 2019.
- 401 [36] Xingzhou Zhang, Yifan Wang, Sidi Lu, Liangkai Liu, Lanyu Xu, and Weisong Shi. Openei: An  
 402 open framework for edge intelligence. *CoRR*, abs/1906.01864, 2019.

## 403 Checklist

- 404 1. For all authors...
- 405 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
 406 contributions and scope? [Yes] See also Section 3, where key contributions were listed.
- 407 (b) Did you describe the limitations of your work? [Yes] See Section 5.3.
- 408 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 409 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
 410 them? [Yes] We read the ethics review guidelines and confirm that content of our paper  
 411 is not related to any ethics issue, since it is focused on the optimization of computations.
- 412 2. If you are including theoretical results...
- 413 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 414 (b) Did you include complete proofs of all theoretical results? [N/A]
- 415 3. If you ran experiments...
- 416 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
 417 mental results (either in the supplemental material or as a URL)? [Yes] See supplemen-  
 418 tal materials.
- 419 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
 420 were chosen)? [N/A] In our paper we do not consider any training, or fine-tuning. We  
 421 focus on the inference only.
- 422 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
 423 ments multiple times)? [N/A] Since we did not do any training, we did not need to run  
 424 experiment multiple times (inference output is deterministic).
- 425 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
 426 of GPUs, internal cluster, or cloud provider)? [N/A] Since we did not do any training,  
 427 it is not related to our work.
- 428 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 429 (a) If your work uses existing assets, did you cite the creators? [Yes] All used assets are  
 430 cited either as a reference, or by direct URL link in the footnote, or in the body of the  
 431 paper.
- 432 (b) Did you mention the license of the assets? [Yes] The license of the used assets are  
 433 noticed either in the appropriate reference, or direct URL link of asset.
- 434 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
 435 Some code and generated LUTs to reproduce our results are provided in supplemental  
 436 materials.
- 437 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
 438 using/curating? [N/A]
- 439 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
 440 information or offensive content? [N/A]
- 441 5. If you used crowdsourcing or conducted research with human subjects...
- 442 (a) Did you include the full text of instructions given to participants and screenshots, if  
 443 applicable? [N/A]
- 444 (b) Did you describe any potential participant risks, with links to Institutional Review  
 445 Board (IRB) approvals, if applicable? [N/A]
- 446 (c) Did you include the estimated hourly wage paid to participants and the total amount  
 447 spent on participant compensation? [N/A]