
decoupleQ: Towards 2-bit Post-Training Uniform Quantization via decoupling Parameters into Integer and Floating Points

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Quantization emerges as one of the most promising compression technologies for
2 deploying efficient large models in recent years. However, existing quantization
3 schemes suffer from significant accuracy degradation at very low bits, or require
4 some additional computational overhead when deployed, making it difficult to be
5 applied to large-scale applications in industry. In this paper, we propose decoupleQ,
6 achieving a substantial increase in model accuracy, especially at very low bits.
7 decoupleQ abandons the traditional heuristic quantization paradigm and decouples
8 the model parameters into integer and floating-point parts, then transforming the
9 quantization problem into a mathematical constrained optimization problem, which
10 is then solved alternatively by off-the-shelf solution methods. decoupleQ gets rid
11 of any tricks for dealing with outliers, sensitive channels, etc., and focuses only
12 on the basic optimization objective to achieve high model accuracy on extreme
13 low bit quantization. Quantization via decoupleQ is linear and uniform, making
14 it hardware-friendlier than non-uniform counterpart, and enabling the idea to be
15 migrated to high-bit quantization to enhance its robustness.
16 decoupleQ has achieved comparable accuracy as fp16/bf16 for 2-bit quantization of
17 large speech models in our company. The code (including the W2 CUDA kernels)
18 is attached and will be made public.

19 1 Introduction

20 Serving large models (1; 2; 37; 38) in industry is budget-consuming because of the huge computa-
21 tional, IO and storage cost. Model compression (10; 11; 16) has therefore become a necessity to
22 alleviate this pain. Among which, Post-Training Quantization (PTQ) (9; 26) has gained more and
23 more popularity among researchers and engineers because it does not require heavy GPU-hours
24 training with labeled datasets.

25 However, previous quantization schemes remain confined within the traditional heuristic quantization
26 paradigm, e.g., how to deal with outliers (33; 35), how to deal with sensitive channels (6), how
27 to determine the clipping range (29), and so on. These methods have achieved some success, but
28 the quantization at extreme low bit often suffers from significant accuracy degradation, thus failing
29 to meet the launching requirements of industrial practice. There are also some other options to
30 mitigate the accuracy loss. QuIP (4) pushes the accuracy limits of 2-bit quantization and can achieve
31 performance close to fp16/bf16. However, compared to traditional quantization schemes, its inference
32 imposes an additional burden due to the need to multiply two random orthogonal matrices to de-
33 quant the weights. N2UQ (20) fit the real-value distribution with non-uniform grids then quantize
34 them into equidistant output levels. But it need to train to get the input thresholds. SpQR (7)

35 and SqueezeLLM (14) use mixed-precision quantization or non-uniform scheme to safeguard the
 36 important channels, but they need customized hardware support.

37 In order to alleviate the above pains in industry, we proposed decoupleQ, which completely abandons
 38 the traditional heuristic quantization paradigm and instead decouples the model parameters into
 39 integer and floating point parts, then transforming the quantization problem into a mathematical
 40 constrained optimization problem, which is then solved alternatively by off-the-shelf solution methods.
 41 The integer part contains the main weights of the model, and the floating-point part contains scales
 42 and zero points induced via quantization. decoupleQ starts from an abstract objective function and
 43 thus does not need any tricks to deal with the minutiae of traditional quantization paradigm, such as
 44 outlier, salient weights (19), and so on. Quantization via decoupleQ is linear and uniform, making it
 45 hardware-friendlier than non-uniform counterpart, and enabling the idea to be migrated to high-bit
 46 quantization to enhance its robustness.

47 decoupleQ contains two stages: 1. layer-wise minimization, defined in Eq. 1, is used to optimize
 48 the integer part and the floating-point part; 2. block-wise minimization, defined in Eq. 2, is used to
 49 further optimize the floating-point part while freezing the integer part¹.

50 Layer-wise minimization is to minimize the ℓ^2 loss of the outputs between pre- and post-quantization
 51 for a linear layer:

$$\min_{\widetilde{W}} \|X\widetilde{W} - XW_0\|_2^2 \quad (1)$$

52 where $X \in \mathbb{R}^{batch \times d_{in}}$ is the input of this layer, $W_0 \in \mathbb{R}^{d_{in} \times d_{out}}$ is the pre-trained full precision
 53 weight, d_{in} and d_{out} are the input and output dimensions respectively. The objective is to find a
 54 matrix \widetilde{W} with quantized-then-dequantized elements to minimize Eq. 1.

55 Some works (4; 8; 9; 13; 25) started from Eq. 1 and achieved some success, but they still haven't
 56 thought outside the box of traditional quantization. GPTQ series (8; 9) fake-quantize the first element
 57 of W_0 and then update the the remaining elements so as to keep Eq. 1 minimized. This process is
 58 then continued element by element until all elements are fake-quantized. However, on the one hand,
 59 they do not give any indication of how scale and zero point should be calculated, and on the other
 60 hand, the optimization problem formulated for updating the remaining elements is unconstrained
 61 (explained in detail later). decoupleQ models Eq. 1 as a constrained optimization problem, as shown
 62 in Eq. 6. It no longer needs to pay attention to some of the minutiae unique to quantization, such as
 63 outliers, clipping threshold, etc., but abstracts the essence of the problem from a higher level.

64 In the second stage, block-wise minimization is used to further improve the model accuracy:

$$\min \|\widetilde{\text{Block}}(X) - \text{Block}(X)\|_2^2 \quad (2)$$

65 where $\widetilde{\text{Block}}(\cdot)$ is a common transformer block (32) with quantized weights. In this stage, we freeze
 66 the integer part of the weights, and train the scales, zero points and norm layers.

67 decoupleQ implements 2-bit uniform quantization and achieves state-of-the-art accuracy in Llama-
 68 1/2 (30; 31). Like traditional uniform quantization, decoupleQ does not incur additional inference
 69 burden and only requires a linear transformation to convert the quantized weights into floating point
 70 ones.

71 Our main highlights are summarized as follows:

- 72 • **New insight:** We abandoned the traditional quantization paradigm, and no longer need
 73 to focus on some of the minutiae unique to quantization, but abstracts the essence of the
 74 problem from a higher level and transforms it into a constrained optimization problem.
- 75 • **Extreme low-bit:** decoupleQ achieves 2-bit uniform quantization with performance match-
 76 ing fp16/bf16 for industrial applications in the ASR model in our company, and we will also
 77 release the W2A16 CUDA kernel as one of our core contribution.
- 78 • **Extensibility:** As a bonus, if labeled datasets are available, the idea of decoupleQ can be
 79 easily extended to supervised fine-tuning (sft) to further improve model accuracy, or the
 80 adaptation to the downstream sub-tasks.

¹We define the term "layer" as a linear transformation, "block" as a common transformer block containing the multi-head attention, feed forward, and some layer norm.

81 2 Related Works

82 Quantization can be roughly divided into Quantization Aware Training (QAT) (21; 33) and Post-
83 Training Quantization (PTQ) (4; 35). In this paper, we focus on weight-only quantization in PTQ,
84 and we will only summarize a few works that are closely related to our work.

85 PTQ is commonly used for LLM quantization because it does not require a lot of GPU hours of
86 training with labeled datasets. AdaRound (25) and BRECQ (18) start from the rounding operation
87 and explore whether to round up or down is better. SqQR (7) and OWQ (17) use mixed-precision
88 quantization strategy to protect sensitive parameters, while AWQ (19) opts for scaling up the weights
89 of sensitive channels to reduce the loss of quantization of sensitive channels. OmniQuant (29) use
90 gradient decent to optimize for the weight clipping threshold and the rescale factors. In decoupleQ, we
91 abandon patchwork solutions and transform the quantization into a principled traditional optimization
92 problem by decoupling the model parameters into integer and floating-point parts.

93 GPTQ (9) is an influential work, and it quantizes the current weights and then updates the remaining
94 weights to minimize the ℓ^2 loss of the output of the layer between pre- and post-quantization. As we
95 will see later, this update actually approximates much, and GPTQ does not optimize for the scale and
96 zero point reduced by quantization.

97 QALora (36) also decouples model parameters at a certain level and uses labeled datasets to fine-tune
98 the zero points. decoupleQ takes this idea a step further, optimizing the scales, zero points and norm
99 layers with supervised fine-tuning, while freezing the integer weights.

100 3 Methods

101 3.1 Preliminaries

102 For a linear layer with input dimension d_{in} and output dimension d_{out} , quantization maps the weights
103 with high-precision into discrete level, and the previous scheme can be described as follows:

$$104 \quad \widehat{W} = \text{clip}(\lfloor \frac{W_0 - z}{s} \rfloor, \alpha, \beta) \quad (3) \quad \widetilde{W} = \widehat{W} * s + z \quad (4)$$

105 where $W_0 \in \mathbb{R}^{d_{in} \times d_{out}}$ is the pre-trained full precision weights, s and z are the scale and zero point
106 (what we call floating-point part above), $\lfloor \cdot \rfloor$ is the round-to-nearest function, $\widehat{W} \in \mathbb{R}^{d_{in} \times d_{out}}$ is the
107 quantized integer-point matrix (what we call integer part above), \widetilde{W} is the de-quantized floating-point
108 matrix, α and β are the lower and upper bounds of the range of integer representations, respectively.
109 For example, in 2-bit weight only linear quantization scheme, the value of each entry of \widehat{W} is
110 limited to one of $\{-2, -1, 0, 1\}$, and $\alpha = -2, \beta = 1$ in this case. To get the values of \widetilde{W} , previous
111 methods (8; 9) show that layer-wise ℓ^2 loss between the outputs pre- and post-quantization is well
112 related to the model accuracy, i.e., to optimize the following objective function,

$$\arg \min_{\widetilde{W}} \|X\widetilde{W} - XW_0\|_2^2 = \text{tr}\{(\widetilde{W} - W_0)^T H(\widetilde{W} - W_0)\} \quad (5)$$

113 where $X \in \mathbb{R}^{batch \times d_{in}}$ is the input of this linear layer, generated by a small set of calibration dataset,
114 and $H = X^T X$.

115 In the very low-bit quantization regime, the model accuracy can be further improved via finer-grained
116 grouping. This would impose additional overhead on inference. For example, when groupsize = 64,
117 it imposes an average overhead of 0.5 bit per element (FP16/BF16 for scale s and zero point z). The
118 extra overhead is acceptable compared to the model accuracy gain.

119 3.2 decoupleQ

120 When a model is quantized, only the integer part \widehat{W} and the floating-point part (s, z) in Eq. 4 are
121 delivered to the downstream inference engine, and the inference process does not need to know how
122 \widehat{W} and (s, z) are obtained at all. That is, if we can find the values of \widehat{W} and (s, z) to minimize Eq. 5
123 by other methods, then we don't need to use Eq. 3. So, we can decouple the model parameters into
124 integer part \widehat{W} and floating point part (s, z) , which are then optimized alternatively via off-the-shelf

125 solution methods. decoupleQ views the process of solving for \widehat{W} and (s, z) in Eq. 4 as a constrained
 126 optimization problem independent of the previous quantization paradigm! We only need to regard
 127 Eq. 4 as an ordinary affine transformation, in which the value of s can be 0 or even negative.

128 In per-channel quantization, each column of the weight matrix is optimized independently of each
 129 other. For simplicity of notation, we only focus on one column in \widehat{W} later and re-define the notations.
 130 Based on Eq. 5, the optimization problem of decoupleQ in the first stage, layer-wise minimization,
 131 can then be formulated as:

$$\begin{aligned} & \min_{w; s, z} g(w; s, z) \\ & \text{s.t. } \forall i = 1, 2, \dots, d_{in} \\ & \quad w_i - \beta \leq 0 \\ & \quad -w_i + \alpha \leq 0 \\ & \quad w_i \in \mathbb{Z} \end{aligned} \quad (6)$$

132 where the objective function is:

$$g(w; s, z) = \frac{1}{2}(w * s + z - b)^T H(w * s + z - b) \quad (7)$$

133 $w \in \mathbb{R}^{d_{in}}$ is one column of \widehat{W} , $b \in \mathbb{R}^{d_{in}}$ is the corresponding column of W_0 , $s \in \mathbb{R}^{ng}$ is the
 134 scale and $z \in \mathbb{R}^{ng}$ is the zero point, ng is the number of groups when grouping-quantization. The
 135 operations w.r.t (s, z) , i.e., $*s$ and $+z$, need to be broadcasted to each group. In this paradigm,
 136 we have completely abandoned the traditional framework of quantization and instead transformed
 137 quantization into a constrained optimization problem 6, which is then solved to achieve the purpose
 138 of quantization. (s, z) in problem 6 have lost the traditional meaning of scale and zero point, and are
 139 just two optimization variables.

140 Transforming the traditional quantization problem into problem 6 is the soul of decoupleQ! Problem 6
 141 is a quadratic programming problem with an additional non-convex constraints $w_i \in \mathbb{Z}$. Quadratic
 142 programming has been studied for many years and there are now many well-established solution (24;
 143 34). We provide one solution in the next subsection, which may not be efficient or optimal.

144 The core idea of decoupleQ is to decouple the model weights into the integer part w and the
 145 floating-point part (s, z) , with the integer part occupying most of the model's expressive power. The
 146 extensibility of the idea of decoupleQ is that we can freeze the integer part of the entire model, and
 147 use labeled data to train the (s, z) as well as other floating point parameters. The advantage of this is
 148 that on the one hand, it can further improve the accuracy of the model, on the other hand, it can fit
 149 specific downstream sub-tasks while maintaining the generalization ability of the model.

150 3.3 Optimization via Alternative Iteration

151 The problem 6 is not easy to solve because of the non-convex constraint $w_i \in \mathbb{Z}$. After obtaining a
 152 good initialization (explained in detail later), we solve for w and (s, z) alternately and iteratively. In
 153 each round of alternation, the objective function 7 w.r.t (s, z) is an unconstrained quadratic function,
 154 thus (s, z) can be readily determined *analytically*: by differentiating the objective function and
 155 equating the derivative to zero, followed by solving the resultant linear system of equations. While
 156 for w , the problem become problem 8:

$$\begin{aligned} & \min_w g(w; s, z) & \min_{w_i; i>j} g(w; s, z) \\ & \text{s.t. } \forall i = 1, 2, \dots, d_{in} & \text{s.t. } \forall i = j + 1, \dots, d_{in} \\ & \quad w_i - \beta \leq 0 & \quad w_i - \beta \leq 0 \\ & \quad -w_i + \alpha \leq 0 & \quad -w_i + \alpha \leq 0 \\ & \quad w_i \in \mathbb{Z} & \quad w_i \in \mathbb{Z} \end{aligned} \quad (8) \quad (9)$$

158 For problem 8, one solution is to round-and-clip one element of w to be integer in $[\alpha, \beta]$ and then
 159 update the remaining. And then this process is then performed sequentially for all elements. After the
 160 j -th element has been rounded-and-clipped, the objective for the updating then becomes problem 9.

161 problem 9 is also intractable, and we can make two levels of approximation:

$$\begin{aligned}
& \min_{w_i; i > j} g(w; s, z) \\
& \text{s.t. } \forall i = j + 1, \dots, d_{in} \\
& w_i - \beta \leq 0 \\
& -w_i + \alpha \leq 0
\end{aligned} \tag{10}$$

$$\min_{w_i; i > j} g(w; s, z) \tag{11}$$

163 In the first-level approximation 10, only the non-convex constraint $w_i \in \mathbb{Z}$ is discarded, while in the
164 second-level approximation 11, both the non-convex constraint $w_i \in \mathbb{Z}$ and the convex constraint
165 $w_i \in [\alpha, \beta]$ are discarded. Intuitively, problem 11 is much simpler to solve than problem 10, but
166 solving problem 10 will lead to a better convergence of the primary objective(6) than solving
167 problem 11. GPTQ (9) provides an efficient analytical solution for problem 11, which we will
168 directly utilize in our experiments. (GPTQ updates the remaining elements by considering only the
169 second-level approximation 11 and ignoring the constrain $w_i \in [\alpha, \beta]$ in the first (10), which is what
170 we mentioned in the introduction, that the update of GPTQ is unconstrained.) As for problem 10,
171 there are many mature solutions in the field of convex optimization, such as active-set method,
172 projected gradient descent (PGD), projected coordinate descent and so on (3). We choose PGD
173 because its parallelization is much better than the other two methods. In the experimental part, we
174 will compare the final accuracy of the model via between solving the first level (10) and the second
175 level 11 approximation on small models, while on large models (e.g. lager than 7 billion parameters),
176 we have to choose the second level 11 approximation because the intolerable runtime of solving the
177 first (10). The algorithm is shown in Alg. 1 and Alg. 2.

Algorithm 1: Alternative Iteration to solve problem 6.

Input: predefined iteration number N .

Result: w^*, s^*, z^*

```

1 Initialize  $t = 1, w_0, s_0, z_0$ ;
2 while  $t \leq N$  do
3   Freeze  $(s_{t-1}, z_{t-1})$ , and optimize
    $g(w; s_{t-1}, z_{t-1})$  to obtain an
178 approximate solution  $w_t$  via
   solving 8 via 2;
4   Freeze  $w_t$ , and solve the
   unconstraint quadratic equation
    $g(w_t; s, z)$  to obtain an analytic
   solution for  $(s_t, z_t)$ ;
5    $t = t + 1$ 
6 end
7  $w^* = w_N; s^* = s_N; z^* = z_N$ 

```

Algorithm 2: Approximate solution of 8

Input: predefined iteration number K, M , and the frozen (s, z) .

Result: w^*

```

1 if Approximaton (10) is used then
2   Ignoring the constraint  $w_i \in \mathbb{Z}$  in Eq. 8, and
   train Eq. 8 with  $M$  iterations via PGD;
3 Initialize  $j = 1$ ;
4 for  $j = 1 \rightarrow d_{in}$  do
5   round and clip the  $j$ -th element of  $w$ , then
   keep the first  $j$  elements frozen, and
   update the remainings via PGD to
   optimize 10 with  $K$  iterations or until
   converged, or via the method in GPTQ to
   optimize 11.
6 end
7  $w^* = w$ 

```

179 3.4 Initialization of w and (s, z)

180 Since the values of w are discrete, a good ini-
181 tialization is very important in order to obtain a
182 more accurate solution to the original problem 6
183 with a faster convergence. Intuitively, the func-
184 tion $g(w; s, z)$ contains the term $w * s$, which
185 means that the scales of the initial values of w
186 and s have to be reasonably distributed. For ex-
187 ample, in the extreme case when the initial value
188 of (s, z) have a very large scale, the first iter-
189 ation will make most of the entries of w strictly
190 0, which will make the iteration crash. We start
191 by initializing (s, z) . We can use grid search to
192 solve the Eq. 12 for the initial value of (s, z) . In
193 Eq. 12, p is a single number, may be different
194 for different columns of W_0 , b_{min} and b_{max} are the minimum and maximum value of b respectively.
195 This step is the same as the previous post-training quantization (19) process. Once the grid search is

$$\begin{aligned}
& \min_p \frac{1}{2} (w * s + z - b)^T H (w * s + z - b) \\
& \text{s.t.} \\
& w = \text{clip}(\lfloor \frac{b - z}{s} \rfloor, \alpha, \beta) \\
& s = \frac{p * (b_{max} - b_{bmin})}{\beta - \alpha} \\
& z = p * b_{min} - s * \alpha
\end{aligned} \tag{12}$$

196 finished, we no longer need to concern ourselves with the (s, z) inside the $\lfloor \cdot \rfloor$ function. The point of
 197 this step is simply to find an initial value for (s, z) for the optimization problem 6.

198 When solving problem 8 via the first-level approximation (10), before entering the for-loop in Alg. 2,
 199 we ignore the constraint $w_i \in \mathbb{Z}$ in problem 8 and optimize it via projected gradient decent with M
 200 iterations. The purpose of this is to allow the first-level approximation to converge in a small number
 201 of iterations, i.e., a small K .

202 3.5 Block-wise minimization

203 After solving problem 6, we obtain a solution for the layer-wise minimization stage and a reasonable
 204 model accuracy. But minimizing the ℓ^2 loss at the layer level does not necessarily lead to the
 205 minimizing the ℓ^2 loss at the block level. We found that the model accuracy can be further improved
 206 via optimization 2. BRECQ (18) also shows that block-reconstruction results in a better model
 207 accuracy than layer-reconstruction. In this stage, we freeze the integer part \widehat{W} in the whole block and
 208 fine-tuning (s, z) and the parameters in norm layer with J epochs.

209 4 Experiments

210 In this section, we describe in detail the experimental results of our method in comparison with other
 211 methods. Unless otherwise stated, all the experiments are conducted on a single A100-SXM-80GB,
 212 and the default experimental setting is as follows:

213 **ResNet:** 10240 images in the training dataloader are used as calibration data, with the standard
 214 augmentation in Pytorch official code (27), and the pretrained full precision checkpoints are from
 215 Torchvision (22). $N = 4, M = 50$ (N and M is defined in `refalg1` and `refalg2`). All the convolution
 216 layers and fully-connected layers are quantized into W2 without groups.

217 **Llama-1/2:** 128 2048-token segments from C4 (28) are used as calibration data. We choose C4
 218 as calibration dataset instead of WikiText2 (23) to be consistent with GPTQ. If the block-wise
 219 minimization is used, we use Adam optimizer (15) to finetune the (s, z) and the parameters in norm
 220 layer with $J = 4$ epochs. The learning rate is $1e-5$, weight decay is $1e-6$.

221 4.1 Private Experiments

222 We applied decoupleQ to our company’s two
 223 Automatic Speech Recognition models(ASR)
 224 (corresponding to task A and task B). Each of
 225 the models contain an encoder and an LLM de-
 226 coder. The input of the models is a speech se-
 227 quence and some prompt, and the output is the
 228 corresponding text. We quantize the LLM de-
 229 coder to W2A16g64. The decoders of the two
 230 models contain 40 transformer blocks with 13
 231 billion parameters and 32 transformer blocks
 232 with 7 billion parameters, respectively. Word Er-
 233 ror Rate (WER) is used as metric to measure the
 234 accuracy of the models (less is better). In this
 235 experiments, we use about 8 millions of speech
 236 tokens as calibration dataset, and train 3 epoch
 237 in each block-wise minimization process. When
 238 an input batch contains sequences of varying
 239 lengths, we use a mask to make sure that the
 240 padding part is not involved in the computation of H and the loss of Eq. 2. In task B, once the whole
 241 model is quantized, we also fine-tune all the (s, z) and layer norm in the LLM with labeled dataset,
 242 while freezing all the integer part \widehat{W} , with 8 A100-SXM-80GB GPUs. The accuracy is shown in
 243 Tab. 1, and the CUDA kernel latency is shown in Fig. 1. The W2A16 CUDA kernel is attached and
 244 will be merged into the NVIDIA repo as one of our core contribution.

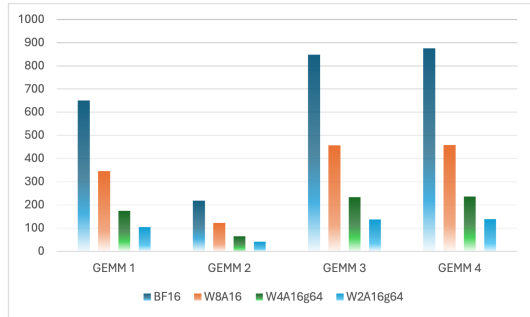


Figure 1: The latency (in $1e-6$ seconds) of the four GEMMs in transformer block on L4 GPU, (The three GEMMs for query, key and value are concatenated into GEMM 1), with $hidden_dim = 5120, batch_size = 4$.

Table 1: The results of our two ASR models. The models are quantized into W2A16g64. runtime for the quantization process is measured in hours. There are two sub-domains in task B, and we report the WER of both.

	Task A		Task B		
	BF16	decoupleQ	BF16	decoupleQ	decoupleQ+sft
WER	6.68	6.70	(5.86, 11.43)	(5.87, 11.56)	(5.77, 11.43)
runtime	-	25	-	32	32+5

Table 2: Comparison of decoupleQ with other methods. In decoupleQ, we only use the first stage, layer-wise minimization. All the models are quantized into W2A16 without groups. In decoupleQ+sft, we train the (s, z) and norm layers for one epoch, using the regular labeled dataset containing 1.2 million images.

method	res18-69.76%			res50-76.13%		
	2bit	3bit	4bit	2bit	3bit	4bit
GPTQ	-	67.88	69.37	-	74.87	75.71
OBQ	64.04	68.69	69.56	70.71	75.24	75.72
BRECQ	64.70	68.47	69.37	72.41	75.32	75.88
decoupleQ	64.15	68.65	69.58	71.34	75.24	76.00
decoupleQ+sft	65.45	68.94	69.71	72.65	75.61	75.97

245 4.2 Public Comparison

246 As a first comparison, we compare decoupleQ with other methods on ImageNet (5) with ResNet (12),
 247 which are standard benchmarks and are efficient to implement. Most importantly, its Top-1 is a strong
 248 indicator of model accuracy. Tab. 2 shows the results of decoupleQ and others. The results other than
 249 decoupleQ are copied from GPTQ (9) and OBQ (8).

250 Tab. 3 shows the results on Llama. In this experiment, we have to choose the second level approxima-
 251 tion(11) because the intolerable runtime of solving the first(10). For a fair comparison, the calibration
 252 dataset contains 128 samples, although a larger calibration dataset will result in stronger results.
 253 we can see that decoupleQ outperforms others almost in all settings, although we use a weaker
 254 approximation(11) to save time. As for the hype-parameters, we choose $\{N = 4, J = 4\}$.

255 4.3 Ablation studies

256 4.3.1 the two approximations

257 The soul of decoupleQ is problem 6, but when solving problem 6, we have to take some approxima-
 258 tions(10 or 11). Obviously, solving approximation 10 will be much more time consuming than solving
 259 approximation 11. But if solving approximation 10 yields better results, the time cost may be worth
 260 it. We first evaluate these two approximations from the perspective of model accuracy. In practice,
 261 we don't have to wait for approximation 10 to fully converge when we solve it via projected gradient
 262 decent, and only need to iterate some steps to get a sub-optimal solution. In Alg. 2, the for-loop takes
 263 up the majority of the runtime. So, we first study the influence of the number of iterations K (defined
 264 in the for-loop) on the final accuracy of the model. Fig. 2 shows the Top-1 accuracy of ResNet-18 on
 265 ImageNet w.r.t the number of iterations K . First of all, in the blue line, we use only the layer-wise
 266 minimization of decooupleQ to quantize the model. After the quantization is finished, in the red line,
 267 we use the labeled dataset with the common 1.2 millions images to fine-tune all the (s, z) and norm
 268 layers for one epoch, with the integer part being frozen. In this step, we use SGD optimizer with
 269 learning rate $1e-6$, weight decaying rate $1e-4$ to train for only one epoch. Fig. 2 clearly indicates the
 270 following conclusions: 1. As the number of iterations K increases, the model accuracy increases
 271 almost monotonically; 2. When $K > 4$, model accuracy via the first approximation(10) is better than
 272 via the second(11). This is to be expected, since the second approximation(11) drops the constraint
 273 $\alpha \leq w_i \leq \beta$, leading to a looser approximation; 3. By the supervised fine-tuning (sft), the model
 274 accuracy is further improved. The same experimental phenomenon also occurs on the ResNet-50
 275 model, which we do not show here.

Table 3: The results of PPL of wikitext-2 on Llama-1/2. We also report the runtime (measured in hours) for the W2 quantization via decoupleQ in the gray background row. The results other than decoupleQ are copied from OmniQuant (29). All the results of decoupleQ use the approximation 11.

Llama		1-7B	1-13B	1-30B	1-65B	2-7B	2-13B	2-70B
FP16		5.68	5.09	4.10	3.53	5.47	4.88	3.31
W2A16	GPTQ	2.1e3	5.5e3	499.75	55.91	7.7e3	2.1e3	77.95
	OmniQuant	15.47	13.21	8.71	7.58	37.37	17.21	7.81
	decoupleQ	9.49	7.86	6.37	5.59	9.74	13.03	5.23
	runtime	2.5	4.8	12.7	27.6	2.5	4.5	33.4
W2A16g128	GPTQ	44.01	15.60	10.92	9.51	36.77	28.14	-
	OmniQuant	8.90	7.34	6.59	5.65	9.62	7.56	6.11
	decoupleQ	8.65	7.25	6.04	5.19	8.79	7.44	4.96
	runtime	3.7	7.7	24.3	55.0	3.7	7.9	70.6
W2A16g64	GPTQ	22.10	10.06	8.54	8.31	20.85	22.44	-
	OmniQuant	8.90	7.34	6.59	5.65	9.62	7.56	6.11
	decoupleQ	8.18	6.96	5.81	5.07	8.41	6.98	5.34
	runtime	4.3	8.9	27.9	64.5	4.4	9.0	98.2
W3A16	GPTQ	8.06	6.76	5.84	5.06	8.37	6.44	4.82
	AWQ	11.88	7.45	10.07	5.21	24.00	10.45	-
	OmniQuant	6.49	5.68	4.74	4.04	6.58	5.58	3.92
	decoupleQ	6.38	5.60	4.67	6.05	6.22	5.72	3.84
W4A16	GPTQ	6.13	5.40	4.48	3.83	5.83	5.13	3.58
	AWQ	6.08	5.34	4.39	3.76	6.15	5.12	-
	OmniQuant	5.86	5.21	4.25	3.71	5.74	5.02	3.47
	decoupleQ	5.85	5.21	4.24	3.67	5.70	5.06	3.45

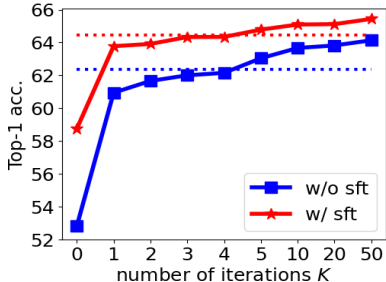


Figure 2: The Top-1 accuracy of ResNet-18 on ImageNet. Solid and dashed lines are for approximation 10 and 11 respectively.

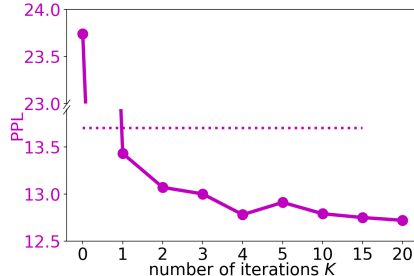


Figure 3: The PPL of Llama-7B on Wiki-Text2. Solid and dashed lines are for approximation 10 and 11 respectively.

276 In the experiment shown in 3, we randomly select 512 2048-token segments from C4 (28). We chose
 277 512 segments here instead of the common 128 in order to reduce the effect of overfitting and thus
 278 compare the two approximations more objectively. In this experiment, we take $N = 2$, and quantize
 279 Llama-7B into W2A16 without groups, and only the layer-wise minimization is used to exclude the
 280 interference of other factors. The PPL decrease almost monotonically as the number of iterations K
 281 increases. It shows that, when $K > 1$, solving approximation 10 yields better model accuracy than
 282 approximation 11.

283 However, when block-wise minimization is introduced in addition to the experiment in 3, the situation
 284 becomes a little more elusive. The results are shown in 4. The model’s best PPL is where $K = 1$,
 285 and then fluctuates within a range as K continues to increase. But all PPLs are inferior to when
 286 the second-level approximation (11) is used. We also plot the loss, defined in 2, of the first block
 287 between pre-and post quantization on the right vertical axis. As K increases, the loss decreases
 288 strictly monotonically, and when $K > 2$, the loss falls below the case when the approximation 11 is
 289 used. This suggests that the correlation between PPL and loss is perhaps weak, and we will investigate
 290 this in the future.

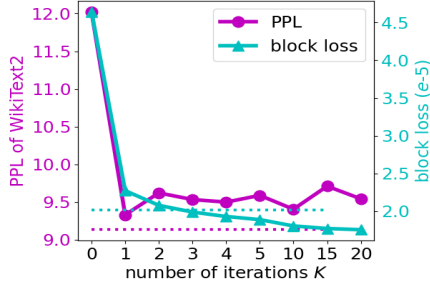


Figure 4: The PPL of Llama-7B on WikiText2 and the loss of the first block between pre- and post-quantization. Solid and dashed lines are for approximation 10 and 11 respectively.

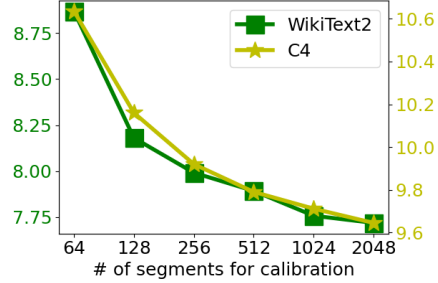


Figure 5: The perplexity of Llama-7B on WikiText2 and C4 dataset w.r.t the number of segments as calibration datasets. The model is quantized into W2A16g64.

291 4.3.2 the size of calibration dataset

292 The solution of problem 6 is dependent on H and thus on the the calibration dataset, as does Eq. 2.
 293 Fig. 5 shows the relationship between dataset size and PPL. In this experiment, Llama-7B is quantized
 294 into W2A16g64. We use the second-level approximation (11) to save time, and $\{N = 4, J = 4\}$. For
 295 runtime reference, when the number of segments is 128/2048, the experiment took 4.3/19.5 hours.

296 4.3.3 the necessity of block-wise minimization

297 Tab. 4 shows that block-wise minimiza-
 298 tion(2) can further improve the model accu-
 299 racy. In this experiment, we choose $N = 4$
 300 and the approximation 11 for the layer-wise
 301 minimization, and $J = 4$ if block-wise min-
 302 imization is used.

Table 4: The perplexity of Llama on WikiText2 with and without the block-wise minimization. All the models are quantized into W2A16.

Llama	1-7B	1-13B	1-30B	2-7B	2-13B
w/o	13.66	9.68	7.35	14.66	12.93
w	9.49	7.86	6.37	9.74	13.03

303 5 Conclusion and Discussion

304 decoupleQ decouples the model parameters into the integer part and a floating point part, and
 305 then optimizes them alternately. This optimization process contains two stages. In the layer-wise
 306 minimization, we transform the quantization problem into the purely mathematical constrained
 307 optimization problem refdecoupleQ; while in the block-wise minimization, we freeze the integer part
 308 and then finetune the floating point part.

309 The risks of decoupleQ include the following: 1. How much the minimization of the ℓ^2 loss of
 310 the layer’s or block’s output correlates with the accuracy of the model; 2. decoupleQ is prone to
 311 overfitting the calibration dataset; 3. The runtime of the quantization process is longer than others.

312 For the first risk, we find experimentally that the correlation between Top-1 and the loss is strong in
 313 the Imagenet classification task; however, the correlation between PPL and the loss is slightly weaker
 314 in LLM. This could be mainly because of an inherent bias between the loss and the accuracy of the
 315 model, or because PPL is not a good indicator of the accuracy of LLM, or for other reasons. For
 316 the second risk, when H in Eq. 7 is an underdetermined matrix, the risk of overfitting rises sharply.
 317 In this case, the possibility of H being underdetermined can be reduced either by enhancing the
 318 diagonal element values of H or by increasing the amount of calibration data. In our practice, we
 319 found that the accuracy of quantization models can rise monotonically with the increase of the size of
 320 the calibration dataset, especially in W2 quantization, but the runtime of quantization rise as well. In
 321 addition, due to time constraints, we do not provide a wealth of public comparisons. However, we
 322 believe that the novelty of a method may outweigh the number of experiments.

323 The idea of decoupleQ is helpful for the adaptation of large model to downstream sub-task. We can
 324 quantize a large foundation model via decoupleQ, then freeze the integer part of the model, and
 325 finetune the floating-point part with labeled dataset from downstream sub-task. Tab. 1 and Tab. 2
 326 show that the model accuracy can be further improved by end-to-end supervised learning.

References

- [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [2] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- [3] Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- [4] Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher De Sa. Quip: 2-bit quantization of large language models with guarantees. *arXiv preprint arXiv:2307.13304*, 2023.
- [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [6] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm.int8(): 8-bit matrix multiplication for transformers at scale. *arXiv preprint arXiv:2208.07339*, 2022.
- [7] Tim Dettmers, Ruslan Svirschevski, Vage Egiazarian, Denis Kuznedelev, Elias Frantar, Saleh Ashkboos, Alexander Borzunov, Torsten Hoefler, and Dan Alistarh. Spqr: A sparse-quantized representation for near-lossless llm weight compression. *arXiv preprint arXiv:2306.03078*, 2023.
- [8] Elias Frantar and Dan Alistarh. Optimal brain compression: A framework for accurate post-training quantization and pruning. *Advances in Neural Information Processing Systems*, 35:4475–4488, 2022.
- [9] Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Optq: Accurate quantization for generative pre-trained transformers. In *The Eleventh International Conference on Learning Representations*, 2022.
- [10] Yi Guo, Yiqian He, Xiaoyang Li, Haotong Qin, Van Tung Pham, Yang Zhang, and Shouda Liu. Rdimkd: Generic distillation paradigm by dimensionality reduction. *arXiv preprint arXiv:2312.08700*, 2023.
- [11] Yi Guo, Huan Yuan, Jianchao Tan, Zhangyang Wang, Sen Yang, and Ji Liu. Gdp: Stabilized neural network pruning via gates with differentiable polarization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5239–5250, 2021.
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [13] Itay Hubara, Yury Nahshan, Yair Hanani, Ron Banner, and Daniel Soudry. Accurate post training quantization with small calibration sets. In *International Conference on Machine Learning*, pages 4466–4475. PMLR, 2021.
- [14] Sehoon Kim, Coleman Hooper, Amir Gholami, Zhen Dong, Xiuyu Li, Sheng Shen, Michael W Mahoney, and Kurt Keutzer. Squeezellm: Dense-and-sparse quantization. *arXiv preprint arXiv:2306.07629*, 2023.
- [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [16] Raghuraman Krishnamoorthi. Quantizing deep convolutional networks for efficient inference: A whitepaper. *arXiv preprint arXiv:1806.08342*, 2018.
- [17] Changhun Lee, Jungyu Jin, Taesu Kim, Hyungjun Kim, and Eunhyeok Park. Owq: Lessons learned from activation outliers for weight quantization in large language models. *arXiv preprint arXiv:2306.02272*, 2023.
- [18] Yuhang Li, Ruihao Gong, Xu Tan, Yang Yang, Peng Hu, Qi Zhang, Fengwei Yu, Wei Wang, and Shi Gu. Breqq: Pushing the limit of post-training quantization by block reconstruction. *arXiv preprint arXiv:2102.05426*, 2021.
- [19] Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Xingyu Dang, and Song Han. Awq: Activation-aware weight quantization for llm compression and acceleration. *arXiv preprint arXiv:2306.00978*, 2023.
- [20] Zechun Liu, Kwang-Ting Cheng, Dong Huang, Eric P Xing, and Zhiqiang Shen. Nonuniform-to-uniform quantization: Towards accurate quantization via generalized straight-through estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4942–4952, 2022.
- [21] Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. Llm-qat: Data-free quantization aware training for large language models. *arXiv preprint arXiv:2305.17888*, 2023.
- [22] Sébastien Marcel and Yann Rodriguez. Torchvision the machine-vision package of torch. In *Proceedings of the 18th ACM international conference on Multimedia*, pages 1485–1488, 2010.
- [23] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*, 2016.
- [24] Katta G Murty and Feng-Tien Yu. *Linear complementarity, linear and nonlinear programming*, volume 3. Heldermann Berlin, 1988.
- [25] Markus Nagel, Rana Ali Amjad, Mart Van Baalen, Christos Louizos, and Tijmen Blankevoort. Up or down? adaptive rounding for post-training quantization. In *International Conference on Machine Learning*, pages 7197–7206. PMLR, 2020.
- [26] Yury Nahshan, Brian Chmiel, Chaim Baskin, Evgenii Zheltonozhskii, Ron Banner, Alex M Bronstein, and Avi Mendelson. Loss aware post-training quantization. *Machine Learning*, 110(11-12):3245–3262, 2021.
- [27] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep

- 394 learning library. *Advances in neural information processing systems*, 32, 2019.
- 395 [28] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou,
396 Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer.
397 *Journal of machine learning research*, 21(140):1–67, 2020.
- 398 [29] Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang, Peng
399 Gao, Yu Qiao, and Ping Luo. Omniquant: Omnidirectionally calibrated quantization for large language
400 models. *arXiv preprint arXiv:2308.13137*, 2023.
- 401 [30] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix,
402 Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation
403 language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 404 [31] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
405 Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and
406 fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 407 [32] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
408 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*,
409 30, 2017.
- 410 [33] Xiuying Wei, Yunchen Zhang, Yuhang Li, Xiangguo Zhang, Ruihao Gong, Jinyang Guo, and Xianglong
411 Liu. Outlier suppression+: Accurate quantization of large language models by equivalent and optimal
412 shifting and scaling. *arXiv preprint arXiv:2304.09145*, 2023.
- 413 [34] Stephen J Wright. *Numerical optimization*. 2006.
- 414 [35] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant:
415 Accurate and efficient post-training quantization for large language models. In *International Conference
416 on Machine Learning*, pages 38087–38099. PMLR, 2023.
- 417 [36] Yuhui Xu, Lingxi Xie, Xiaotao Gu, Xin Chen, Heng Chang, Hengheng Zhang, Zhensu Chen, Xiaopeng
418 Zhang, and Qi Tian. Qa-lora: Quantization-aware low-rank adaptation of large language models. *arXiv
419 preprint arXiv:2309.14717*, 2023.
- 420 [37] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher
421 Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models.
422 *arXiv preprint arXiv:2205.01068*, 2022.
- 423 [38] Yu Zhang, Wei Han, James Qin, Yongqiang Wang, Ankur Bapna, Zhehuai Chen, Nanxin Chen, Bo Li,
424 Vera Axelrod, Gary Wang, et al. Google usm: Scaling automatic speech recognition beyond 100 languages.
425 *arXiv preprint arXiv:2303.01037*, 2023.

426 NeurIPS Paper Checklist

427 1. Claims

428 Question: Do the main claims made in the abstract and introduction accurately reflect the
429 paper's contributions and scope?

430 Answer: [Yes]

431 Justification: Our claims and justification include:

- 432 **a.** Our results are higher than others in very low bit (2-bit) quantization.(This is
433 justified in Tab. 3.);
- 434 **b.** decoupleQ has achieved comparable accuracy as fp16/bf16 for 2-bit quantiza-
435 tion of large speech models in our company. (This is justified in Tab. 1, and the
436 W2 CUDA kernel used in our company are attached.);
- 437 **c.** decoupleQ gets rid of any tricks for dealing with outliers, sensitive channels,
438 etc. (This is justified in the Problem 6, we do not use any tricks, such as scaling
439 factor (19; 29), mixed-precision quantization (6), etc., to deal with outliers and
440 sensitive channels.)

441 Guidelines:

- 442 • The answer NA means that the abstract and introduction do not include the claims
443 made in the paper.
- 444 • The abstract and/or introduction should clearly state the claims made, including the
445 contributions made in the paper and important assumptions and limitations. A No or
446 NA answer to this question will not be perceived well by the reviewers.
- 447 • The claims made should match theoretical and experimental results, and reflect how
448 much the results can be expected to generalize to other settings.
- 449 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
450 are not attained by the paper.

451 2. Limitations

452 Question: Does the paper discuss the limitations of the work performed by the authors?

453 Answer: [Yes]

454 Justification: The paper has discussed the three limitations of decoupleQ in the last section,
455 **Conclusion and Discussion**, and the risk overall that we did not provide as many public
456 comparison experiments as other work due to time constraints.

457 Guidelines:

- 458 • The answer NA means that the paper has no limitation while the answer No means that
459 the paper has limitations, but those are not discussed in the paper.
- 460 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 461 • The paper should point out any strong assumptions and how robust the results are to
462 violations of these assumptions (e.g., independence assumptions, noiseless settings,
463 model well-specification, asymptotic approximations only holding locally). The authors
464 should reflect on how these assumptions might be violated in practice and what the
465 implications would be.
- 466 • The authors should reflect on the scope of the claims made, e.g., if the approach was
467 only tested on a few datasets or with a few runs. In general, empirical results often
468 depend on implicit assumptions, which should be articulated.
- 469 • The authors should reflect on the factors that influence the performance of the approach.
470 For example, a facial recognition algorithm may perform poorly when image resolution
471 is low or images are taken in low lighting. Or a speech-to-text system might not be
472 used reliably to provide closed captions for online lectures because it fails to handle
473 technical jargon.
- 474 • The authors should discuss the computational efficiency of the proposed algorithms
475 and how they scale with dataset size.
- 476 • If applicable, the authors should discuss possible limitations of their approach to
477 address problems of privacy and fairness.

478 • While the authors might fear that complete honesty about limitations might be used by
479 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
480 limitations that aren't acknowledged in the paper. The authors should use their best
481 judgment and recognize that individual actions in favor of transparency play an impor-
482 tant role in developing norms that preserve the integrity of the community. Reviewers
483 will be specifically instructed to not penalize honesty concerning limitations.

484 3. Theory Assumptions and Proofs

485 Question: For each theoretical result, does the paper provide the full set of assumptions and
486 a complete (and correct) proof?

487 Answer:[NA]

488 Justification: This paper does not include theoretical results.

489 Guidelines:

- 490 • The answer NA means that the paper does not include theoretical results.
- 491 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
492 referenced.
- 493 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 494 • The proofs can either appear in the main paper or the supplemental material, but if
495 they appear in the supplemental material, the authors are encouraged to provide a short
496 proof sketch to provide intuition.
- 497 • Inversely, any informal proof provided in the core of the paper should be complemented
498 by formal proofs provided in appendix or supplemental material.
- 499 • Theorems and Lemmas that the proof relies upon should be properly referenced.

500 4. Experimental Result Reproducibility

501 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
502 perimental results of the paper to the extent that it affects the main claims and/or conclusions
503 of the paper (regardless of whether the code and data are provided or not)?

504 Answer: [Yes]

505 Justification: At the beginning of the section **Experiments**, we provide details of the experi-
506 mental parameters; specifically for each experiment, we also provide the key experimental
507 parameters.

508 Guidelines:

- 509 • The answer NA means that the paper does not include experiments.
- 510 • If the paper includes experiments, a No answer to this question will not be perceived
511 well by the reviewers: Making the paper reproducible is important, regardless of
512 whether the code and data are provided or not.
- 513 • If the contribution is a dataset and/or model, the authors should describe the steps taken
514 to make their results reproducible or verifiable.
- 515 • Depending on the contribution, reproducibility can be accomplished in various ways.
516 For example, if the contribution is a novel architecture, describing the architecture fully
517 might suffice, or if the contribution is a specific model and empirical evaluation, it may
518 be necessary to either make it possible for others to replicate the model with the same
519 dataset, or provide access to the model. In general, releasing code and data is often
520 one good way to accomplish this, but reproducibility can also be provided via detailed
521 instructions for how to replicate the results, access to a hosted model (e.g., in the case
522 of a large language model), releasing of a model checkpoint, or other means that are
523 appropriate to the research performed.
- 524 • While NeurIPS does not require releasing code, the conference does require all submis-
525 sions to provide some reasonable avenue for reproducibility, which may depend on the
526 nature of the contribution. For example
 - 527 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
528 to reproduce that algorithm.
 - 529 (b) If the contribution is primarily a new model architecture, the paper should describe
530 the architecture clearly and fully.

- 531 (c) If the contribution is a new model (e.g., a large language model), then there should
532 either be a way to access this model for reproducing the results or a way to reproduce
533 the model (e.g., with an open-source dataset or instructions for how to construct
534 the dataset).
- 535 (d) We recognize that reproducibility may be tricky in some cases, in which case
536 authors are welcome to describe the particular way they provide for reproducibility.
537 In the case of closed-source models, it may be that access to the model is limited in
538 some way (e.g., to registered users), but it should be possible for other researchers
539 to have some path to reproducing or verifying the results.

540 5. Open access to data and code

541 Question: Does the paper provide open access to the data and code, with sufficient instruc-
542 tions to faithfully reproduce the main experimental results, as described in supplemental
543 material?

544 Answer: [Yes]

545 Justification: The code (including W2 CUDA kernels) is attached in supplementary material,
546 and can reproduce the results in the public experiments.

547 Guidelines:

- 548 • The answer NA means that paper does not include experiments requiring code.
- 549 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
550 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 551 • While we encourage the release of code and data, we understand that this might not be
552 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
553 including code, unless this is central to the contribution (e.g., for a new open-source
554 benchmark).
- 555 • The instructions should contain the exact command and environment needed to run to
556 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
557 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 558 • The authors should provide instructions on data access and preparation, including how
559 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 560 • The authors should provide scripts to reproduce all experimental results for the new
561 proposed method and baselines. If only a subset of experiments are reproducible, they
562 should state which ones are omitted from the script and why.
- 563 • At submission time, to preserve anonymity, the authors should release anonymized
564 versions (if applicable).
- 565 • Providing as much information as possible in supplemental material (appended to the
566 paper) is recommended, but including URLs to data and code is permitted.

567 6. Experimental Setting/Details

568 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
569 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
570 results?

571 Answer: [Yes]

572 Justification: At the beginning of the section **Experiments**, we provide details of the experi-
573 mental parameters; specifically for each experiment, we also provide the key experimental
574 parameters. The code is attached in supplementary material and will be made public.

575 Guidelines:

- 576 • The answer NA means that the paper does not include experiments.
- 577 • The experimental setting should be presented in the core of the paper to a level of detail
578 that is necessary to appreciate the results and make sense of them.
- 579 • The full details can be provided either with the code, in appendix, or as supplemental
580 material.

581 7. Experiment Statistical Significance

582 Question: Does the paper report error bars suitably and correctly defined or other appropriate
583 information about the statistical significance of the experiments?

584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634

Answer: [No]

Justification: The cost of the experiment is high.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We have reported that most of the experiments are conducted in one single A100-SXM-80GB, except for the sft process. And we also reported the time of execution.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

Answer: [Yes]

Justification: The research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: This is a work for accelerating the inference of deep models, where the social impact is determined by the function of the model, not by how the inference is accelerated.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This work does not release models or datasets.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The paper has cited all related works, and included the relevant license.

687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [\[Yes\]](#)

Justification: We provide the source code, and a readme and license file are alongside.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [\[NA\]](#)

Justification: the paper does not involve crowdsourcing nor research with human subjects

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [\[NA\]](#)

739
740
741
742
743
744
745
746
747
748
749
750

Justification: the paper does not involve crowdsourcing nor research with human subjects

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.