# ATTENTION-AWARE POST-TRAINING QUANTIZATION WITHOUT BACKPROPAGATION

Anonymous authors

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

023

024 025

048

051

052

Paper under double-blind review

# ABSTRACT

Quantization offers a promising solution for deploying large-scale language models (LLMs) on resource-constrained devices. However, early quantization methods, developed for smaller networks like ResNet, rely on gradient-based optimization, which becomes impractical for hyper-scale LLMs with billions of parameters. While recently proposed backpropagation-free post-training quantization (PTQ) methods alleviate this issue, their performance is limited by a lack of interlayer dependency consideration. In this paper, we introduce a novel PTQ algorithm that incorporates inter-layer dependencies without relying on backpropagation. The key innovation is the development of attention-aware Hessian matrices that capture inter-layer interactions within the attention module. Extensive experiments demonstrate that our approach significantly outperforms conventional PTQ methods, particularly at low bit-widths.

1 INTRODUCTION

The explosive growth in complexity (parameters) of large-scale language models (LLMs) based on Transformers (Touvron et al., 2023; Zhang et al., 2022) has resulted in a proportional increase in computational costs, which has prompted an urgent need for efficient model processing and compression strategies. Quantization has emerged as a pivotal solution in this context, and it serves as an essential step in deploying AI models on resource-constrained devices that primarily support fixed-point arithmetic. By reducing precision, the memory bandwidth requirements can be alleviated, and the significant parallelism of quantized models can be SIMDified using highly efficient vector processing units, such as neural processing units (NPUs).

Two main categories of quantization approaches have been proposed to preserve the performance of original full-precision models: quantization-aware training (QAT) and post-training quantization (PTQ). Although QAT can potentially outperform PTQ, its practicality diminishes considerably when handling hyper-scale LLMs featuring billions of parameters. Consequently, recent quantization efforts have been directed toward PTQ.

Although classic PTQ methods have successfully quantized small-scale models (Nagel et al., 2020; Li et al., 2021), they rely on time-consuming gradient-based optimization, so their efficacy decreases when the complexity of LLMs increases. Accordingly, backpropagation-free PTQ methods have been developed for LLMs (Frantar et al., 2023; Xiao et al., 2023; Jeon et al., 2023b); however, their performance is somewhat limited owing to the lack of consideration of inter-layer dependencies. Recent studies have attempted to consider inter-layer dependencies (Shao et al., 2023; Ma et al., 2024), but they still rely on time-consuming gradient-based optimizations.

In this paper, we propose a novel quantization algorithm that considers inter-layer dependencies without relying on backpropagation. Our primary contributions can be summarized as follows:

• We propose a novel PTQ algorithm called BOA<sup>1</sup>. To avoid time-consuming gradientbased optimization, we adopt the Hessian-based strategy introduced by (Frantar & Alistarh, 2022). The primary contribution is to approximate the Hessian more accurately by exploiting the attention reconstruction error, not the layer-wise reconstruction error, to capture inter-layer dependencies within the attention module (**Section 3.2**).

<sup>&</sup>lt;sup>1</sup>BOA: <u>Backpropagation-free optimization for Attention-aware PTQ</u>

• While the proposed Hessian facilitates the consideration of inter-layer dependencies, it requires a large amount of memory and high computational cost. Therefore, we incorporate several techniques to mitigate the computational overhead, including Hessian relaxation, efficient computation of inverse Hessians, and head-wise simultaneous quantization (Section 3.3).

• We evaluate BOA via extensive experiments on publicly available LLMs. Our results demonstrate that BOA outperforms conventional LLM PTQ methods by a significant margin, particularly for low-bit precision (*e.g.*, INT2) (Section 4).

## 2 RELATED WORKS

054

056

059

060

061

062 063

064 065

066

067

068

069

071

072

073

074

075 076 077 When calibration data are available, PTQ primarily aims to minimize the increase in task loss incurred by quantization. Consider a neural network parameterized by weights **W**. Provided that the network is trained to convergence, the problem of quantizing **W** to minimize task loss degradation can be formulated as (LeCun et al., 1989)

$$\min_{\Delta \mathbf{w}} \Delta \mathbf{w}^T \cdot \mathbf{H}^{(\mathbf{w})} \cdot \Delta \mathbf{w}, \tag{1}$$

where  $\mathbf{H}^{(\mathbf{w})}$  is the Hessian related to the flattened weight  $\mathbf{w}$  and  $\Delta \mathbf{w}$  is a weight perturbation caused by the quantization. Owing to the infeasibility in computing and storing the exact Hessian  $\mathbf{H}^{(\mathbf{w})}$ , many studies have assumed independence between layers, which relaxes (1) into the following layerwise reconstruction problem (Nagel et al., 2020):

$$\min_{\mathbf{W}^{(\ell)}} \left\| \Delta \mathbf{W}^{(\ell)} \mathbf{X}^{(\ell-1)} \right\|_{F}^{2}, \tag{2}$$

where  $\mathbf{X}^{(\ell-1)}$  is the input to the  $\ell$ -th layer parameterized by  $\mathbf{W}^{(\ell)}$ .

079 To solve (2), early efforts aimed to optimize the weight-rounding mechanism (Nagel et al., 2020; Hubara et al., 2021; Li et al., 2021; Jeon et al., 2022; 2023a). Instead of allocating the nearest quanti-081 zation bin, these studies attempted to assign quantized values that minimize the reconstruction error. 082 In (Nagel et al., 2020), a backpropagation-based optimization method, called AdaRound, has been 083 proposed. This algorithm has been extended to BRECQ where the block-wise reconstruction error 084 has been used instead of the layer-wise reconstruction error to consider the inter-layer dependencies 085 within a certain network block (e.g., Transformer block), which leads to an enhanced low-bit quantization performance (Li et al., 2021). Although AdaRound and BRECQ have successfully quantized small-sized models such as ResNet (He et al., 2016), they are heavily dependent on time-consuming 087 gradient-based optimizations. This renders their application to LLMs with billions of parameters 880 challenging. Consequently, recent efforts have shifted towards the development of cost-effective 089 quantization methods for LLMs. 090

091 These efforts can be classified into two orthogonal categories: 1) Hessian-based methods that opti-092 mize a weight-rounding mechanism without relying on backpropagation (e.g., GPTQ (Frantar et al., 2023)) and 2) equivalent transformation (ET)-based methods that transform a model to be robust 093 to quantization, thereby enhancing the performance of the naive rounding-to-nearest quantization 094 (e.g., SmoothQuant (Xiao et al., 2023), AWQ (Lin et al., 2024), Z-FOLD (Jeon et al., 2023b), Out-095 lier Suppression+ (OS+) (Wei et al., 2023), OmniQuant (Shao et al., 2023), AffineQuant (Ma et al., 096 2024)). Among them, GPTQ has emerged as one of the most efficient quantization methods; GPTQ is capable of quantizing models with over 30 billion parameters in just a few GPU hours, with negli-098 gible performance degradation at INT4 precision. Moreover, GPTQ can be integrated with ET-based methods to enhance their performance (Lin et al., 2024; Jeon et al., 2023b).<sup>2</sup> Similar to GPTQ, we 100 optimize a weight-rounding mechanism based on Hessian without relying on backpropagation. The 101 primary difference is that we pursue the preservation of the attention output after the quantization 102 while GPTQ aims to preserve each layer output and thus cannot consider inter-layer dependencies 103 within the attention module.

Other algorithms exploiting different quantization strategies have also been proposed. For example, SpQR (Dettmers et al., 2023), SqueezeLLM (Kim et al., 2023), and OAC (Edalati et al., 2024) tried

106 107

<sup>2</sup>By transforming models with ET-based methods and then optimizing a weight-rounding mechanism with GPTQ (instead of applying the rounding-to-nearest method), the quantization performance can be boosted.

108 to preserve quantization-sensitive weights by assigning a large bit-width or retaining them in full-109 precision. QuIP (Chee et al., 2023) introduced the idea of incoherent processing to suppress outliers 110 within weights. When compared to the standard uniform quantization, these algorithms require ad-111 ditional processing and memory costs in the real inference stage and need some special hardware and 112 dedicated kernels without which accelerating the inference process may not be easy.<sup>3</sup> Furthermore, unlike server-grade GPUs (e.g. NVIDIA A100), on-device NPUs (e.g. Qualcomm Hexagon) lack 113 support for the mixed precision format and such additional processing, and customizing kernels for 114 desired functionalities is very challenging on on-device NPUs. Thus, we exclude these algorithms 115 in our comparison and focus on the more universally supported uniform quantization format. 116

#### 3 METHOD

117 118

119 120

121

122

130

131 132

133

138 139

148

149 150

151

152

161

### 3.1 OVERVIEW OF PROPOSED BOA

Similar to GPTQ, the proposed BOA algorithm quantizes weights by repeating the quantization and weight-update steps; once BOA quantizes one weight, it updates the remaining (not-yet-quantized) 123 weights by exploiting the Hessian-based weight-update formula introduced by GPTQ, compensating 124 for the task loss degradation caused by the quantization. When the q-th weight  $w_q$  is quantized, the 125 weight-update  $\delta w$  is mathematically expressed as

$$\delta \boldsymbol{w} = -\frac{w_q - \mathcal{Q}(w_q)}{[\mathbf{U}]_{q,q}} [\mathbf{U}]_{q,:} \text{ where } \mathbf{U} = \text{Chol}(\mathbf{H}^{-1})^T.$$
(3)

Here, **H** is the Hessian,  $Chol(\cdot)$  denotes a Cholesky decomposition (*i.e.*, **U** is an upper triangular matrix satisfying  $\mathbf{H}^{-1} = \mathbf{U}^T \mathbf{U}$ , and  $\mathcal{Q}$  is a uniform quantization function defined as

$$\mathcal{Q}(x) = s \left( \operatorname{clamp} \left( \left\lfloor \frac{x}{s} \right\rfloor + z, 0, 2^n - 1 \right) - z \right),$$

where s, z, n are the scale, zero-point, and bit-width, respectively, and  $|\cdot|$  represents the round-off. 134

135 The key difference over GPTQ lies in the approximation of the Hessian H. In GPTQ, the layer-wise 136 independence has been assumed to approximate H, which yields the following Hessian equation<sup>4</sup>: 137

$$\mathbf{H}^{(\boldsymbol{w}^{(\ell)})} \approx 2\mathbf{X}^{(\ell-1)}\mathbf{X}^{(\ell-1)^{T}} \otimes \mathbf{I},\tag{4}$$

where  $\mathbf{H}^{(w^{(\ell)})}$  is the Hessian for the  $\ell$ -th layer,  $\otimes$  denotes the Kronecker product operation and  $\mathbf{I}$  is 140 the identity matrix. As the approximated Hessian  $\mathbf{H}^{(\boldsymbol{w}^{(\ell)})}$  relies solely on the input, GPTQ is unable 141 to account for the influence of other layers when compensating for quantization error (see (3)). In 142 other words, GPTQ neglects inter-layer dependencies within the attention module, a crucial aspect 143 of Transformers (Vaswani et al., 2017), which results in somewhat constrained performance for low 144 precision (e.g., INT2) (Jeon et al., 2023b). To overcome this, we develop Hessians that incorporate 145 inter-layer dependencies within the attention module and then use them instead of the conventional 146 Hessian in (4). 147

### 3.2 PROPOSED ATTENTION-AWARE HESSIAN

To consider the inter-layer dependencies within the attention module, we exploit the attention reconstruction error rather than the layer-wise reconstruction error when approximating the Hessian.

For an input sequence  $\mathbf{X} \in \mathbb{R}^{d \times L}$ , the output of the multi-head attention (MHA) is expressed as

$$MHA(\mathbf{X}) = \sum_{h=1}^{H} \mathbf{W}_{\text{out},h} (\mathbf{A}_h \mathbf{V}_h)^T, \ \mathbf{A}_h = \sigma \left(\frac{\mathbf{Q}_h \mathbf{K}_h^T}{\sqrt{d_h}}\right),$$
(5)

<sup>157</sup> <sup>3</sup>In QuIP, the weight W is multiplied by random orthogonal matrices U and V (*i.e.*, W  $\leftarrow$  UWV<sup>T</sup>). 158 While this incoherent processing can suppress outliers within weights, additional processing is needed to re-159 cover quantized weights (*i.e.*,  $\widehat{\mathbf{W}} \leftarrow \mathbf{U}^T \widehat{\mathbf{W}} \mathbf{V}$ ; see Algorithm 2 in (Chee et al., 2023)). Such processing must be done during the inference, incurring additional inference time and memory costs for storing U and V. 160

<sup>&</sup>lt;sup>4</sup>For any  $\mathbf{M}_1$  and  $\mathbf{M}_2$ , the second-order derivative of  $\|\mathbf{M}_1 \Delta \mathbf{W} \mathbf{M}_2\|_F^2$  with respect to  $\Delta \mathbf{w}$  is  $2\mathbf{M}_2 \mathbf{M}_2^T \otimes$  $\mathbf{M}_{1}^{T}\mathbf{M}_{1}$  (see Appendix A for the proof).

179 180

181

182

183 184

185

186

187 188 189

190

199

200 201

202

204 205 206

208 209 210

213

#### Algorithm 1 BOA 163 **Input**: weights $\mathbf{W} \in \mathbb{R}^{d_{row} \times d_{col}}$ and inputs **X** of the Transformer layer 164 1: def BOA(W, X)Initialize quantized output: $\mathbf{Q} \leftarrow \mathbf{0}_{H \times d_{\text{row}}/H \times d_{\text{col}}}$ 2: 166 Initialize (row-wise) quantization errors: $\mathbf{E} \leftarrow \mathbf{0}_{H \times d_{col}}$ 3: 167 $\triangleright$ See Table 1 4: Compute attention-aware Hessians: $\mathbf{H}_h = \mathbf{H}_{\text{col},h} \otimes \mathbf{H}_{\text{row},h}$ 168 Set step size (scale) S: min<sub>S</sub> tr $(\Delta WH_{col,h}\Delta W^T)$ 5: 169 Compute inverse Hessians $\mathbf{H}_{\text{col},h}^{-1}$ and $\mathbf{H}_{\text{row},h}^{-1}$ 6: 170 Compute $\mathbf{U}_{\text{col},h} = \text{Chol}(\mathbf{H}_{\text{col},h}^{-1})^T$ and $\mathbf{U}_{\text{row},h} = \text{Chol}(\mathbf{H}_{\text{row},h}^{-1})^T$ 7: 171 8: for $j = 0, ..., d_{row}/H - 1$ do 172 Construct $\mathbf{W}^{(j)} \in \mathbb{R}^{H \times d_{col}}$ by stacking the *j*-th rows $[\mathbf{W}_h]_{j,:}$ 9: 173 Quantize $\mathbf{W}^{(j)}$ : $(\mathbf{Q}_{:,j,:}, \mathbf{E}) \leftarrow \text{GPTQ}(\mathbf{W}^{(j)}, \mathbf{U}_{\text{col},h}, \mathbf{S})$ 10: ▷ See Appendix D 174 Update remaining rows: $[\mathbf{W}_h]_{j:,:} \leftarrow [\mathbf{W}_h]_{j:,:} - \frac{[\mathbf{U}_{\text{row},h}^T]_{j:,j} \cdot \mathbf{E}_{h,:} \cdot \mathbf{U}_{\text{col},h}}{[\mathbf{U}_{\text{row},h}]_{j,j}}$ 175 11: 176 12: end for 177 Output: quantized weights Q 178

where  $\sigma$  is the row-wise softmax function, H is the number of attention heads,  $d_h$  is the embedding dimension of the h-th attention head,  $[\mathbf{Q}_h|\mathbf{K}_h|\mathbf{V}_h] = \mathbf{X}^T[\mathbf{W}_{Q,h}^T|\mathbf{W}_{K,h}^T|\mathbf{W}_{V,h}^T], \mathbf{W}_{\{Q,K,V\},h} \in$  $\mathbb{R}^{d_h \times d}$ , and  $\mathbf{W}_{\text{out},h} \in \mathbb{R}^{d \times d_h}$ .

**Hessian for**  $\mathbf{W}_{Q,h}$  When  $\mathbf{W}_{Q,h}$  is quantized,  $\mathbf{W}_{\text{out},h}$  and  $\mathbf{V}_h$  in (5) remain unchanged, but the attention weight  $A_h$  changes. Using the first-order Taylor polynomial, the perturbation in  $A_h$  can be approximated as

$$\Delta \mathbf{A}_{h} = \sigma \left( \frac{(\mathbf{Q}_{h} + \Delta \mathbf{Q}_{h})\mathbf{K}_{h}^{T}}{\sqrt{d_{h}}} \right) - \sigma \left( \frac{\mathbf{Q}_{h}\mathbf{K}_{h}^{T}}{\sqrt{d_{h}}} \right) \approx \frac{\Delta \mathbf{Q}_{h}\mathbf{K}_{h}^{T}}{\sqrt{d_{h}}} \mathbf{J}_{\sigma}^{T} = \frac{\mathbf{X}^{T} \Delta \mathbf{W}_{Q,h}^{T}\mathbf{K}_{h}^{T}\mathbf{J}_{\sigma}^{T}}{\sqrt{d_{h}}}, \quad (6)$$

where  $\mathbf{J}_{\sigma}$  is the Jacobian matrix of the softmax function  $\sigma$ . Thus, the attention reconstruction error is expressed as

$$\|\Delta \mathrm{MHA}(\mathbf{X})\|_{F}^{2} = \|\mathbf{W}_{\mathrm{out},h}(\Delta \mathbf{A}_{h}\mathbf{V}_{h})^{T}\|_{F}^{2} = \left\|\frac{\mathbf{W}_{\mathrm{out},h}\mathbf{V}_{h}^{T}\mathbf{J}_{\sigma}\mathbf{K}_{h}}{\sqrt{d_{h}}}\Delta \mathbf{W}_{Q,h}\mathbf{X}\right\|_{F}^{2}.$$
 (7)

Combining this with Footnote 4 yields the following Hessian for  $W_{Q,h}$ :

$$\mathbf{H}^{(\mathbf{w}_{Q,h})} = 2\mathbf{X}\mathbf{X}^{T} \otimes \frac{\mathbf{K}_{h}^{T}\mathbf{J}_{\sigma}^{T}\mathbf{V}_{h}\mathbf{W}_{\text{out},h}^{T}\mathbf{W}_{\text{out},h}\mathbf{V}_{h}^{T}\mathbf{J}_{\sigma}\mathbf{K}_{h}}{d_{h}}.$$
(8)

**Hessian for**  $\mathbf{W}_{K,h}$  When  $\mathbf{W}_{K,h}$  is quantized, the attention weight  $\mathbf{A}_h$  changes as in the quantization of  $\mathbf{W}_{Q,h}$ . By following the steps for (6),  $\mathbf{A}_h$  can be approximated as

$$\Delta \mathbf{A}_{h} \approx \frac{\mathbf{Q}_{h} \Delta \mathbf{K}_{h}^{T}}{\sqrt{d_{h}}} \mathbf{J}_{\sigma}^{T} = \frac{\mathbf{Q}_{h} \Delta \mathbf{W}_{K,h} \mathbf{X} \mathbf{J}_{\sigma}^{T}}{\sqrt{d_{h}}},\tag{9}$$

203 and thus the attention reconstruction error can be expressed as

$$\|\Delta \mathrm{MHA}(\mathbf{X})\|_{F}^{2} = \|\Delta \mathbf{A}_{h} \mathbf{V}_{h} \mathbf{W}_{\mathrm{out},h}^{T}\|_{F}^{2} = \left\|\frac{\mathbf{Q}_{h}}{\sqrt{d_{h}}} \Delta \mathbf{W}_{K,h} \mathbf{X} \mathbf{J}_{\sigma}^{T} \mathbf{V}_{h} \mathbf{W}_{\mathrm{out},h}^{T}\right\|_{F}^{2}.$$
 (10)

Thus, we obtain the following Hessian for  $\mathbf{W}_{K,h}$ : 207

$$\mathbf{H}^{(\mathbf{w}_{K,h})} = 2\mathbf{X}\mathbf{J}_{\sigma}^{T}\mathbf{V}_{h}\mathbf{W}_{\text{out},h}^{T}\mathbf{W}_{\text{out},h}\mathbf{V}_{h}^{T}\mathbf{J}_{\sigma}\mathbf{X}^{T} \otimes \frac{\mathbf{Q}_{h}^{T}\mathbf{Q}_{h}}{d_{h}}.$$
(11)

When quantizing  $\mathbf{W}_{V,h}$ , only  $\mathbf{V}_h$  changes. Thus, the attention reconstruction **Hessian for**  $\mathbf{W}_{V,h}$ 211 error is expressed as 212

$$\|\Delta \operatorname{MHA}(\mathbf{X})\|_{F}^{2} = \|\mathbf{W}_{\operatorname{out},h}(\mathbf{A}_{h}\Delta\mathbf{V}_{h})^{T}\|_{F}^{2} = \|\mathbf{W}_{\operatorname{out},h}\Delta\mathbf{W}_{V,h}\mathbf{X}\mathbf{A}_{h}^{T}\|_{F}^{2},$$

214 which yields the following Hessian for  $\mathbf{W}_{V,h}$ : 215

$$\mathbf{H}^{(\mathbf{w}_{V,h})} = 2\mathbf{X}\mathbf{A}_{h}^{T}\mathbf{A}_{h}\mathbf{X}^{T} \otimes \mathbf{W}_{\text{out},h}^{T}\mathbf{W}_{\text{out},h}.$$
(12)

Table 1: Proposed attention-aware Hessians

	-
Layer	н
$\mathbf{W}_{Q,h}$	$2\mathbf{X}\mathbf{X}^T\otimes\mathbf{K}_h^T\mathbf{K}_h$
$\mathbf{W}_{K,h}$	$2\mathbf{X}\mathbf{X}^T\otimes \mathbf{Q}_h^T\mathbf{Q}_h$
$\mathbf{W}_{V,h}$	$2\mathbf{X}\mathbf{A}_{h}^{T}\mathbf{A}_{h}\mathbf{X}^{T}\otimes\mathbf{W}_{\mathrm{out},h}^{T}\mathbf{W}_{\mathrm{out},h}$
$\mathbf{W}_{\mathrm{out},h}$	$2\mathbf{X}_{\mathrm{out},h}\mathbf{X}_{\mathrm{out},h}^T\otimes\mathbf{I}$
$\mathbf{W}_{\mathrm{fc1}}$	$2\mathbf{X}_{ ext{fc1}}\mathbf{X}_{ ext{fc1}}^T\otimes \mathbf{I}$
$\mathbf{W}_{\mathrm{fc2}}$	$2\mathbf{X}_{ ext{fc2}}\mathbf{X}_{ ext{fc2}}^T\otimes \mathbf{I}$

**Hessian for**  $\mathbf{W}_{\text{out},h}$  When  $\mathbf{W}_{\text{out},h}$  is quantized, the attention reconstruction error is expressed as

$$\|\Delta \operatorname{MHA}(\mathbf{X})\|_{F}^{2} = \|\Delta \mathbf{W}_{\operatorname{out},h}(\mathbf{A}_{h}\mathbf{V}_{h})^{T}\|_{F}^{2}.$$

Thus, the corresponding Hessian is obtained as

$$\mathbf{H}^{(\mathbf{w}_{out,h})} = 2\mathbf{V}_{h}^{T}\mathbf{A}_{h}^{T}\mathbf{A}_{h}\mathbf{V}_{h} \otimes \mathbf{I} = 2\mathbf{X}_{out}\mathbf{X}_{out}^{T} \otimes \mathbf{I},$$
(13)

where  $\mathbf{X}_{\text{out},h} = (\mathbf{A}_h \mathbf{V}_h)^T$  is the input to the out-projection layer.

3.3 EFFICIENT IMPLEMENTATION OF BOA

While inter-layer dependencies within the attention module can be considered by exploiting the proposed Hessians, they are significantly more complex than the conventional Hessian in (4), which may incur high computational costs. For example, computing the proposed Hessians in (8) and (11) would be more expensive than computing the conventional one in (4). In this subsection, we present techniques for mitigating the computational overheads incurred by the proposed attention-aware Hessians.

241

249 250

251

257

260 261

262

269

216

224 225

226 227 228

229 230 231

232 233

234

**Hessian relaxation** The largest overhead related to the computation of the proposed Hessians is the Jacobian matrix  $\mathbf{J}_{\sigma}$  in (8) and (11). For one input sequence, the shape of  $\mathbf{J}_{\sigma}$  is  $H \times L \times L \times L$ , which requires a large amount of memory and high computational cost (more than 400 GB even for the OPT-125M model when L = 2048).

To mitigate such computational overhead, we establish a relaxed Hessian that does not require the computation of  $J_{\sigma}$ . To this end, we build an upper bound for the attention reconstruction error in (7), which will be used as its surrogate:

$$\|\Delta \operatorname{MHA}(\mathbf{X})\|_{F}^{2} \leq \left\|\frac{\mathbf{W}_{\operatorname{out},h}\mathbf{V}_{h}^{T}\mathbf{J}_{\sigma}}{\sqrt{d_{h}}}\right\|_{F}^{2} \cdot \left\|\mathbf{K}_{h}\Delta \mathbf{W}_{Q,h}\mathbf{X}\right\|_{F}^{2}.$$
(14)

Moreover, we note that the term  $\|\mathbf{W}_{out,h}\mathbf{V}_{h}^{T}\mathbf{J}_{\sigma}\|_{F}^{2}$  in (14) is constant and does not affect quantization.<sup>5</sup> Thus, we do not need to consider this term when computing the Hessian. In short, we use the term  $\|\mathbf{K}_{h}\Delta\mathbf{W}_{Q,h}\mathbf{X}\|_{F}^{2}$  as a surrogate of the attention reconstruction error when deriving the Hessian for  $\mathbf{W}_{Q,h}$ , which results in the following relaxed Hessian:

$$\mathbf{H}^{(\mathbf{w}_{Q,h})} = 2\mathbf{X}\mathbf{X}^T \otimes \mathbf{K}_h^T \mathbf{K}_h.$$
(15)

258 259 Similarly, we can establish a relaxed Hessian for  $\mathbf{W}_{K,h}$  as follows:

$$\mathbf{H}^{(\mathbf{w}_{K,h})} = 2\mathbf{X}\mathbf{X}^T \otimes \mathbf{Q}_h^T \mathbf{Q}_h.$$
(16)

In Table 1, we summarize the relaxed Hessians for each layer inside the Transformer block.

Efficient computation of inverse Hessians Owing to the size of the proposed attention-aware Hessians being  $dd_h \times dd_h$ , the complexity of the computation of the inverse Hessian (see (3)) would be  $\mathcal{O}(d^3d_h^3)$  in our approach. This is considerably more expensive than the complexity  $\mathcal{O}(d^3)$  in GPTQ, where the inverse of only the column-wise Hessian  $\mathbf{X}\mathbf{X}^T \in \mathbb{R}^{d \times d}$  (in (4)) is needed (Frantar et al., 2023).

<sup>&</sup>lt;sup>5</sup>The weight-update  $\delta \mathbf{w}$  in (3) is not affected by the constant multiple of **H** because  $[c\mathbf{U}]_{q,:}/[c\mathbf{U}]_{q,q} = [\mathbf{U}]_{q,:}/[\mathbf{U}]_{q,q}$  for any constant *c*.



Figure 1: Illustration of the proposed BOA for the query projection  $W_{Q}$ .

For the efficient computation of the inverse Hessians, we exploit the useful properties of the Kronecker product (see (17a)-(17c) in Appendix A). For simplicity, let  $\mathbf{H} = \mathbf{H}_{col} \otimes \mathbf{H}_{row}$  where  $\mathbf{H}_{col} \in \mathbb{R}^{d \times d}$  and  $\mathbf{H}_{row} \in \mathbb{R}^{d_h \times d_h}$ , then we obtain 

$$\mathbf{H}^{-1} = (\mathbf{H}_{col} \otimes \mathbf{H}_{row})^{-1} = \mathbf{H}_{col}^{-1} \otimes \mathbf{H}_{row}^{-1}$$

This implies that the inverse Hessian  $\mathbf{H}^{-1}$  can be computed by computing  $\mathbf{H}_{col}^{-1}$  and  $\mathbf{H}_{row}^{-1}$  (line 6 in Algorithm 1) whose complexity is  $\mathcal{O}(d^3) + \mathcal{O}(d_h^3)$  (=  $\mathcal{O}(d^3)$ ), not  $\mathcal{O}(d^3d_h^3)$ . Similarly, we can efficiently compute the Cholesky decomposition with the same order of complexity as in GPTQ. Specifically, if  $\mathbf{L}_1 = \text{Chol}(\mathbf{H}_{col}^{-1})$  and  $\mathbf{L}_2 = \text{Chol}(\mathbf{H}_{row}^{-1})$ ,  $\mathbf{H}^{-1}$  can be expressed as 

 $\mathbf{H}^{-1} = \mathbf{L}_1 \mathbf{L}_1^T \otimes \mathbf{L}_2 \mathbf{L}_2^T = (\mathbf{L}_1 \otimes \mathbf{L}_2) (\mathbf{L}_1 \otimes \mathbf{L}_2)^T.$ 

Subsequently, noting that the Kronecker product of lower triangular matrices is also lower triangular, we obtain 

$$\operatorname{Chol}(\mathbf{H}^{-1}) = \mathbf{L}_1 \otimes \mathbf{L}_2 = \operatorname{Chol}(\mathbf{H}_{\operatorname{col}}^{-1}) \otimes \operatorname{Chol}(\mathbf{H}_{\operatorname{row}}^{-1}).$$

Thus, we can obtain  $\text{Chol}(\mathbf{H}^{-1})$  by computing  $\text{Chol}(\mathbf{H}_{\text{col}}^{-1})$  and  $\text{Chol}(\mathbf{H}_{\text{row}}^{-1})$  (line 7 in Algorithm 1). Consequently, the computational complexity of the Choleksy decomposition would be  $\mathcal{O}(d^3)$ , not  $\mathcal{O}(d^3 d_h^3).$ 

Simultaneous quantization of different heads The conventional Hessian in (4) implies the independence of different rows (see  $H_{row} = I$ ). Whereas, the proposed attention-aware Hessians model the dependency between different rows (e.g.,  $\mathbf{H}_{row} = \mathbf{K}_h^T \mathbf{K}_h$  for the query projection), using which we can compensate for the quantization error of a certain row by updating the other rows. However, in this case, the rows must be quantized sequentially (not simultaneously). For example, the second row can be quantized after being updated to compensate for the quantization error of the first row. This is in contrast to GPTQ where all the rows are quantized simultaneously.

To accelerate the quantization process, we assume independence between different attention heads (see Fig. 1(a)), under which rows related to different heads are independent and can thus be quan-tized together. For a better understanding, we consider the query projection  $\mathbf{W}_Q$  as an example (see Fig. 1(b)). In the quantization step, we stack the *j*-th rows  $[\mathbf{W}_{Q,h}]_{j,:}$  of all different heads, construct-ing the sub-weight matrix  $\mathbf{W}_{Q}^{(j)} \in \mathbb{R}^{H \times d}$  (line 9 in Algorithm 1). Because the rows of  $\mathbf{W}_{Q}^{(j)}$  are mutually independent, all the rows of  $\mathbf{W}_{Q}^{(j)}$  can be quantized simultaneously as in GPTQ (line 10 in Algorithm 1). Following the quantization of j-th rows, we compensate for the quantization error by updating the remaining rows. In this update step, we use the refined weight-update formula (line 11 in Algorithm 1); the detailed derivation for this is provided in Appendix B.

324 To evaluate how much the quantization process can be accelerated by the head-wise simultaneous 325 quantization, we measure the processing times of the proposed method with and without simultane-326 ous quantization. From Table 2, we observe that BOA requires a significantly long processing time 327 without the simultaneous quantization (more than one day for the 2.7B model). This is because all 328 rows need to be quantized sequentially (e.g., 2560 rows are quantized sequentially for OPT-2.7b) and thus the massive compute capabilities of modern GPUs cannot be utilized properly. As evident 329 from Table 2, by applying the head-wise simultaneous quantization, we can achieve a significant 330 reduction in the processing time. 331

 Table 2: Processing time of BOA with and without head-wise simultaneous quantization

Head-wise	Model Size (OPT)							
Simultaneous Quantization	125M	350M	1.3B	2.7B				
Х	49.22 min	181.7 min	712.7 min	24.48 hr				
0	5.099 min	13.93 min	31.64 min	1.101 hr				

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETUP

344 To evaluate the performance of the proposed BOA, we quantize publicly available language models including OPT (Zhang et al., 2022), BLOOM (Scao et al., 2022), and LLaMA (Touvron et al., 345 2023)). As in (Frantar et al., 2023; Jeon et al., 2023b; Chee et al., 2023), we quantize only weights 346 and retain activations with full precision because activations do not pose a significant bottleneck 347 for the inference of LLMs (Frantar et al., 2023; Kim et al., 2023). As a calibration dataset, we use 348 128 random 2048 token segments from the C4 dataset (Raffel et al., 2020). Thus, we do not use 349 any task-specific data for quantization. We evaluate the performance of the quantized models using 350 benchmark datasets (e.g., WikiText-2 (Merity et al., 2016), C4 (Raffel et al., 2020), and PTB (Marcus 351 et al., 1993)) and zero-shot tasks. All experiments were conducted using a single NVIDIA A100 352 GPU (80 GB). 353

When determining a quantization order in BOA, the heuristic introduced by GPTQ can be employed; the column/row corresponding to the largest diag( $\mathbf{H}_{col}$ )/diag( $\mathbf{H}_{row}$ ) (*i.e.*, the most quantizationsensitive column/row) is first quantized for better compensation. Empirically, we observed that this heuristic could occasionally enhance the performance, yet at other times, it may result in inferior performance. We conduct experiments with and without this heuristic and report the better results.

359

332

340

341 342

343

## 4.2 COMPARISON WITH GPTQ

360 361

We compare the proposed BOA with GPTQ (Frantar et al., 2023), which is our primary baseline. 362 For both algorithms, we set per-channel quantization parameters (*i.e.*, scale and zero-point) to min-363 imize the layer-wise reconstruction error (line 5 in Algorithm 1). We note that in GPTQ, the Min-364 Max-based quantization parameters have been used (Frantar et al., 2023); however, this results in significantly worse quantization performance (Jeon et al., 2023b). While both algorithms aim to op-366 timize the weight-rounding mechanism and can be combined with existing ET-based methods such 367 as SmoothQuant (Xiao et al., 2023), AWQ (Lin et al., 2024), and Z-FOLD (Jeon et al., 2023b), we 368 do not perform an equivalent transform in this experiment to solely compare the weight-rounding optimization performance. The results of integration with ET-based methods are presented in Sec-369 tion 4.3. 370

First, we compare the perplexity (PPL) performances of BOA and GPTQ (see Table 3 and Tables 7 and 8 in Appendix C.1). The performance of the rounding-to-nearest (RTN) method (which naively assigns the nearest quantized value) is also included for comparison, as in (Frantar et al., 2023). While RTN collapses for low bit-widths, BOA and GPTQ exhibit reasonable PPL, even for INT2 quantization. This is because BOA and GPTQ aim to minimize the task loss degradation, not the weight quantization error  $\Delta W$ . Evidently, the proposed BOA outperforms GPTQ for all models. In particular, the performance gap is significant for low bit-width (*i.e.*, INT2) and small-sized models suited for resource-limited devices (*e.g.*, mobile devices).

379

418

419

420

421

422

423 424 425

380					(4	i) OF 1				
381		Dataset	Method	1251	M 350N	1 1.3B	2.7B	6.7B	13B	30B
382			RTN	5.5e	3 2.8e4	1.1e5	9.5e3	2.8e4	1.9e5	1.7e5
383		WikiText-2	GPTQ	232.	8 98.65	66.76	37.44	24.74	18.97	13.12
384			BOA	141.	6 57.40	48.71	26.20	22.71	18.76	12.15
385			RTN	4.3e	3 2.8e4	1.1e4	6.8e3	1.8e4	1.2e5	1.7e5
386		PTB	GPTQ	384.	8 135.9	9 112.0	64.59	42.36	26.95	20.25
387			BOA	199.	2 90.87	78.73	40.76	33.77	25.34	18.52
388			RTN	3.7e	3 1.6e4	7.7e3	7.7e3	1.4e4	9.7e4	5.8e4
389		C4	GPTQ	178.	6 71.89	64.11	33.94	24.86	20.08	14.45
390			BOA	118.	1 54.07	48.92	26.57	23.03	19.22	13.84
391				(	b) BLOO	M and L	LaMA			
392							-			
393		Dataset	Method	5601	M 11B	BLOOM 1 7B	/I 3R	7 1 B	L 13B	LaMA 30B
394			DTN	7.0-	5 0.0-5	2.5.5	1.4-5	2.1.5	5.7-4	300
395		WikiText-2	GPTO	7.8e	3 43.93	3.5e5 36.48	1.4e5 29.25	2.165	5.7e4 12.67	+ 2.7e4 7 8.844
396		Wiki loke 2	BOA	52.0	9 38.16	<b>30.76</b>	24.25	17.54	11.56	5 7.993
397			RTN	7 4e	5 1.1ef	5 2.5e5	1.2e5	2.2e5	8 1e4	1 3 3e4
398		PTB	GPTQ	142.	.6 176.4	4 95.32	67.48	43.73	20.55	5 14.64
399			BOA	113.	.0 139.1	6 <b>9.98</b>	53.10	35.97	18.49	) 13.24
400			RTN	1.4e	6 2.1e6	5 2.7e5	9.2e4	1.3e5	5.9e4	1 2.8e4
401		C4	GPTQ	57.3	1 43.48	38.69	30.97	23.52	14.24	4 11.78
402			BOA	52.1	2 39.03	3 33.71	27.26	21.22	13.34	4 10.53
403		* INT3/INT4 c	quantization	results a	are provide	d in Appe	ndix C.1 du	ue to the p	page limita	tion.
404										
405										
406										
407	Table	e 4: INT2 zer	o-shot tas	sk perf	ormance	(accura	cy ↑) of t	he prop	osed BO	A and GP
408		Model	Size M	lethod	ARC-c	ARC-e	HellaSw	ag MN	MLU AV	/erage
409			G	рто	22 53	35.61	34.03	22	93 2	878
410			<sup>1.3B</sup> <b>B</b>	oA	22.53	38.72	36.00	23	3.12 <b>3</b>	0.09
411		-	G	PTO	24.40	38 /17	37.87	23	3 0/ 3	0.95
412			2.7B B	oA	25.51	42.89	43.68	23	3.14 <b>3</b>	<b>3.81</b>
413		-	G	PTO	25.60	12.85	13 20	2/	100 3	3.96
414		OPT	6.7B B	0A	26.62	44.91	44.52	2- 24	1.33 <b>3</b>	5.10
415		-	G	ΡΤΟ	26.62	44 15	50.00	2/	150 3	6 36
416			<sup>13B</sup> <b>B</b>	0A	27.47	47.39	54.42	25	5.21 <b>3</b>	8.62
417		-		DTO	01.57	52.00	(0.55			2 (0

Table 3: INT2 quantization performance (PPL  $\downarrow$ ) of the proposed BOA and the conventional GPTQ. (a) OPT

ΡTQ.

52.99

53.24

58.71

59.01

62.84

63.47

60.55

62.58

57.48

59.73

65.09

66.31

25.27

26.41

23.53

23.90

31.16

33.11

42.60

43.43

42.97

44.11

49.05

50.19

426	Next, we compare the zero-shot performances of BOA and GPTQ (see Table 4). To this end, we
427	measure the accuracy of quantized models for several tasks and then average the results. We note
428	that the zero-shot setting is maintained in our experiments because we do not use task-specific data
429	for quantization. As evident, the proposed BOA outperforms GPTQ for all models. The key factor
430	leading to such an outstanding performance is that we consider inter-layer dependencies within the
431	attention module by targeting attention-wise reconstruction. This is in contrast to GPTQ, where
	layers are assumed to be independent.

GPTQ

BOA

GPTQ

BOA

GPTQ

BOA

30B

13B

30B

LLaMA

31.57

31.48

32.17

33.79

37.12

37.88

-									
Equivalent	Dataset	Method			Mod	el Size (C	OPT)		
Transformation	Dataset	wictilou	125M	350M	1.3B	2.7B	6.7B	13B	30B
	WiltiToyt 2	GPTQ	229.4	N/A	39.88	27.31	20.03	15.32	13.55
	WIKITCAT-2	BOA	151.2	N/A	31.62	24.45	18.55	14.29	12.49
Same eth Oracent	DTD	GPTQ	292.3	N/A	64.17	44.75	32.01	22.03	19.26
SmoothQuant	FID	BOA	223.9	N/A	58.17	38.87	27.85	19.79	17.97
	C4	GPTQ	151.4	N/A	38.13	26.80	21.22	16.19	14.42
		BOA	130.4	N/A	34.20	24.95	20.92	15.33	13.90
		GPTQ	156.0	102.5	33.97	27.10	18.07	16.29	13.24
	wiki lext-2	BOA	107.9	54.72	29.38	23.96	17.18	15.14	12.41
Z-Fold	DTD	GPTQ	206.9	130.7	53.80	46.08	26.79	23.73	19.27
	PID	BOA	166.1	82.27	49.18	39.45	24.94	22.86	18.11
	<u> </u>	GPTQ	108.8	71.37	31.67	25.98	19.79	17.21	14.13
		BOA	86.07	49.39	28.65	24.19	19.01	16.17	13.67

Table 5: INT2 performance (PPL  $\downarrow$ ) of BOA integrated with existing ET-based methods.

\* SmoothQuant does not support OPT-350M where the post-LayerNorm architecture has been used.

### 448 449 450

432

# 4.3 INTEGRATION WITH EQUIVALENT TRANSFORM-BASED METHODS

As mentioned, the performance of the proposed BOA can be enhanced by combining BOA with existing ET-based methods (*i.e.*, transforming models with ET-based methods first and then applying BOA for optimizing the weight-rounding mechanism). To verify this, we evaluate the performance of BOA integrated with ET-based methods. Among various algorithms, we use SmoothQuant (Xiao et al., 2023) and Z-FOLD (Jeon et al., 2023b) in our integration because they efficiently find out an equivalent transform without time-consuming gradient-based optimization.<sup>6</sup>

Table 5 and Table 9 (see Appendix C.2) summarize the PPL performances of the proposed BOA
combined with SmoothQuant and Z-FOLD. For comparison, we also summarize the integration
results for the conventional GPTQ. Overall, the performance of BOA indeed improves when combined with ET-based methods. We emphasize that the performance gap between the proposed BOA
and GPTQ still remains significant for INT2 quantization. A similar behavior can be observed in the
zero-shot results (see Table 10 in Appendix C.2); the performance is boosted by applying ET-based
methods, and BOA outperforms GPTQ for all models regardless of the ET-based method.

403 464 465

## 4.4 COMPARISON WITH PRIOR ARTS

466 We compare the proposed BOA with OmniQuant (Shao et al., 2023) and AffineQuant (Ma et al., 467 2024), recently proposed algorithms that learn an attention-aware equivalent transform via back-468 propagation (see Table 6 and Table 11 in Appendix C.3 for PPL results and see Table 12 in Ap-469 pendix C.3 for zero-shot results). In our comparison, we do not include AWQ (Lin et al., 2024) 470 and OS+ (Wei et al., 2023) because they perform worse than OmniQuant and AffineQuant (Shao 471 et al., 2023; Ma et al., 2024). Mixed quantization algorithms (e.g., SpQR (Dettmers et al., 2023), 472 SqueezeLLM (Kim et al., 2023), and OAC (Edalati et al., 2024)) and algorithms that require addi-473 tional processing in the real inference stage (e.g., QuIP (Chee et al., 2023)) are also not included 474 because they require some dedicated kernels for acceleration which may not be supported by ondevice NPUs such as Qualcomm Hexagon. 475

As evident, BOA itself outperforms existing algorithms in almost all cases, even though BOA does not rely on time-consuming gradient-based optimization and thus facilitates fast quantization (see Table 13 in Appendix C.4). Furthermore, when combined with SmoothQuant or Z-FOLD, the performance gap between BOA and OmniQuant/AffineQuant is significant, which demonstrates the efficacy of the proposed BOA. We observe that OmniQuant and AffineQuant sometimes diverge or collapse (*i.e.*, PPL is larger than 10<sup>3</sup>) for INT2 quantization. In fact, to supplement the INT2 quantization performance, group-wise quantization parameters have been additionally used in OmniQuant

<sup>&</sup>lt;sup>6</sup>While SmoothQuant has been proposed in the context of weight-activation quantization, the smoothing factor  $(\mathbf{s}_j = \max(|\mathbf{X}_j|)^{\alpha} / \max(|\mathbf{W}_j|)^{1-\alpha})$  used for the equivalent transformation can also be used for weightonly quantization by setting  $\alpha = 0$ .

Dataset	Method	125M	1.3B	2.7B	6.7B	13B	30E
	OmniQuant	NaN	NaN	NaN	2.3e4	4.5e5	3.86
	AffineQuant	174.5	NaN	42.26	26.25	38.89	5.6
WikiText-2	BOA	141.6	48.71	26.20	22.71	18.76	12.
	BOA + SmoothQuant	151.2	31.62	24.45	18.55	14.29	12.4
	BOA + Z-Fold	107.9	29.38	23.96	17.18	15.14	12.4
	OmniQuant	NaN	NaN	NaN	5.0e4	3.7e5	2.9
	AffineQuant	254.2	NaN	55.58	37.36	50.10	3.1
PTB	BOA	199.2	78.73	40.76	33.77	25.34	18.
	BOA + SmoothQuant	223.9	58.17	38.87	27.85	19.79	17.
	BOA + Z-Fold	166.1	49.18	39.45	24.94	22.86	18.
	OmniQuant	NaN	NaN	NaN	3.0e4	2.0e5	2.1
	AffineQuant	107.0	NaN	34.45	25.11	31.50	3.3
C4	BOA	118.1	48.92	26.57	23.03	19.22	13.
	BOA + SmoothQuant	130.4	34.20	24.95	20.92	15.33	13.
	BOA + Z-FOLD	86.07	28.65	24.19	19.01	16.17	13.

Table 6: INT2 performance (PPL  $\downarrow$ ) of BOA and existing approaches.

\* 'NaN' means that loss diverges in the quantization process.

and AffineQuant, but group-wise parameters result in additional memory costs and processing time in the real inference step (Shen et al., 2023).

### 508 4.5 Comparison of Time and Memory Costs

We compare the processing time and memory costs of BOA and conventional algorithms (see Ta-ble 13 in Appendix C.4). We observe that the processing time of BOA is shorter than those required by existing attention-aware algorithms (*i.e.*, OmniQuant and AffineQuant), yet BOA achieves sig-nificantly better performance (see Tables 6 and 11), which demonstrates the efficacy of the proposed method. We also observe that BOA requires longer processing time and larger memory than those required by GPTQ. This is because GPTQ quantizes all the rows of the weight matrix simultaneously using only layer input. In contrast, BOA sequentially quantizes sub-weight matrices using outputs of other layers as well as layer input (see Fig. 1(b)) to consider the inter-layer dependencies within the attention module, which eventually leads to the better quantization performance than GPTQ. 

Clearly, there is a trade-off between quantization speed / memory cost and accuracy. In real situa-tions, when one needs to preserve the performance of the original model by considering inter-layer dependencies within the attention module, the proposed BOA would be an intriguing solution. Even when the memory resource is limited, BOA can be used with some relaxation. Specifically, we note that the large memory cost of BOA for hyper-scale LLMs (e.g., 13B and 30B) is attributable to the row-wise Hessian for the value projection  $(\mathbf{X}\mathbf{A}_h^T\mathbf{A}_h\mathbf{X}^T \text{ in (12)})$  whose shape is  $H \times d \times d$ . In memory-limited cases, we can mitigate the memory cost of BOA by considering inter-layer depen-dencies only for query and key projections and applying the standard Hessian ( $\mathbf{X}\mathbf{X}^{T}$  in (4)) for the value projection. Indeed, when applying the proposed Hessians only for query and key projections, BOA requires almost same amount of memory as GPTQ, yet still exhibiting better performance (see Table 14 in Appendix C.4). For more discussion on time and memory costs of the proposed BOA, see Appendix C.4.

## 5 CONCLUSION

In this paper, we proposed a novel PTQ algorithm called BOA. To consider the inter-layer dependencies within the attention module while circumventing time-consuming gradient-based optimization, we approximated the Hessian matrices by exploiting the attention reconstruction error. Furthermore, to mitigate the computational overhead incurred by the proposed attention-aware Hessians, we incorporated several techniques, such as Hessian relaxation, efficient computation of inverse Hessians, and head-wise simultaneous quantization. Finally, through extensive experiments, we demonstrated the efficacy of the proposed BOA algorithm.

#### 540 REFERENCES 541

565

583

- Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher De Sa. QuIP: 2-bit quantization 542 of large language models with guarantees. In Thirty-seventh Conference on Neural Information 543 Processing Systems, 2023. 544
- Tim Dettmers, Ruslan Svirschevski, Vage Egiazarian, Denis Kuznedelev, Elias Frantar, Saleh Ashk-546 boos, Alexander Borzunov, Torsten Hoefler, and Dan Alistarh. SpQR: A sparse-quantized repre-547 sentation for near-lossless llm weight compression. arXiv:2306.03078, 2023.
- 548 Ali Edalati, Alireza Ghaffari, Masoud Asgharian, Lu Hou, Boxing Chen, and Vahid Partovi Nia. 549 OAC: Output-adaptive calibration for accurate post-training quantization. arXiv:2405.15025, 550 2024. 551
- 552 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, 553 2022. 554
- 555 Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. OPTQ: Accurate quantization 556 for generative pre-trained Transformers. In The Eleventh International Conference on Learning Representations, 2023. 558
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-559 nition. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition 560 (CVPR), pp. 770–778, 2016. 561
- 562 Itay Hubara, Yury Nahshan, Yair Hanani, Ron Banner, and Daniel Soudry. Accurate post training 563 quantization with small calibration sets. In International Conference on Machine Learning, pp. 564 4466-4475. PMLR, 2021.
- Yongkweon Jeon, Chungman Lee, Eulrang Cho, and Yeonju Ro. Mr. BiQ: Post-training non-uniform 566 quantization based on minimizing the reconstruction error. In Proceedings of the IEEE/CVF 567 Conference on Computer Vision and Pattern Recognition, pp. 12329–12338, 2022. 568
- 569 Yongkweon Jeon, Chungman Lee, and Ho-young Kim. GENIE: show me the data for quantization. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 570 12064–12073, 2023a. 571
- 572 Yongkweon Jeon, Chungman Lee, Kyungphil Park, and Ho-young Kim. A frustratingly easy post-573 training quantization scheme for LLMs. In Proceedings of the 2023 Conference on Empirical 574 Methods in Natural Language Processing, pp. 14446–14461, 2023b. 575
- Sehoon Kim, Coleman Hooper, Amir Gholami, Zhen Dong, Xiuyu Li, Sheng Shen, Michael W 576 Mahoney, and Kurt Keutzer. SqueezeLLM: Dense-and-sparse quantization. arXiv:2306.07629, 577 2023. 578
- 579 Yann LeCun, John S Denker, Sara A Solla, Richard E Howard, and Lawrence D Jackel. Optimal 580 brain damage. In Advances in Neural Information Processing Systems (NIPS), volume 2, pp. 581 598-605, 1989.
- 582 Yuhang Li, Ruihao Gong, Xu Tan, Yang Yang, Peng Hu, Qi Zhang, Fengwei Yu, Wei Wang, and Shi Gu. BRECQ: Pushing the limit of post-training quantization by block reconstruction. In 584 International Conference on Learning Representations (ICLR), 2021.
- Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan 586 Xiao, Xingyu Dang, Chuang Gan, and Song Han. AWQ: Activation-aware weight quantization 587 for LLM compression and acceleration. In MLSys, 2024. 588
- 589 Yuexiao Ma, Huixia Li, Xiawu Zheng, Feng Ling, Xuefeng Xiao, Rui Wang, Shilei Wen, Fei Chao, 590 and Rongrong Ji. AffineQuant: Affine transformation quantization for large language models. 591 arXiv:2403.12544, 2024. 592
- Mitchell Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a large annotated corpus of english: The penn treebank. 1993.

- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. *arXiv:1609.07843*, 2016.
- 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 (ICML)*, pp. 7197–7206, 2020.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
   Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text
   transformer. *Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman
  Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. BLOOM: A 176Bparameter open-access multilingual language model. *arXiv:2211.05100*, 2022.
- Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang,
   Peng Gao, Yu Qiao, and Ping Luo. OmniQuant: Omnidirectionally calibrated quantization for
   large language models. *arXiv:2308.13137*, 2023.
- Haihao Shen, Hanwen Chang, Bo Dong, Yu Luo, and Hengyu Meng. Efficient llm inference on CPUs. *arXiv:2311.00502*, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
  Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. LLaMA: Open and
  efficient foundation language models. *arXiv:2302.13971*, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Kiuying Wei, Yunchen Zhang, Yuhang Li, Xiangguo Zhang, Ruihao Gong, Jinyang Guo, and Xianglong Liu. Outlier suppression+: Accurate quantization of large language models by equivalent and effective shifting and scaling. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 1648–1665, 2023.
- Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. SmoothQuant:
   Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*, pp. 38087–38099. PMLR, 2023.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christo pher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. OPT: Open pre-trained Transformer
   language models. *arXiv:2205.01068*, 2022.

# A PROOF OF FOOTNOTE 4

In our proof, we use the following useful properties of the Kronecker product:

$$\operatorname{vec}\left(\mathbf{M}_{1}\mathbf{M}_{2}\mathbf{M}_{3}\right) = \left(\mathbf{M}_{3}^{T} \otimes \mathbf{M}_{1}\right)\operatorname{vec}(\mathbf{M}_{2}),\tag{17a}$$

$$\left(\mathbf{M}_1 \otimes \mathbf{M}_2\right)^T = \mathbf{M}_1^T \otimes \mathbf{M}_2^T, \tag{17b}$$

$$\left(\mathbf{M}_{1}\otimes\mathbf{M}_{2}\right)\left(\mathbf{M}_{3}\otimes\mathbf{M}_{4}\right)=\mathbf{M}_{1}\mathbf{M}_{3}\otimes\mathbf{M}_{2}\mathbf{M}_{4},\tag{17c}$$

where  $vec(\cdot)$  denotes the vectorization operation.

Using (17a), we have

$$\left\|\mathbf{M}_{1}\Delta\mathbf{W}\mathbf{M}_{2}\right\|_{F}^{2} = \left\|\left(\mathbf{M}_{2}^{T}\otimes\mathbf{M}_{1}\right)\Delta\mathbf{w}\right\|_{2}^{2} = \Delta\mathbf{w}^{T}\left(\mathbf{M}_{2}^{T}\otimes\mathbf{M}_{1}\right)^{T}\left(\mathbf{M}_{2}^{T}\otimes\mathbf{M}_{1}\right)\Delta\mathbf{w},$$

where  $\Delta \mathbf{w} = \text{vec}(\Delta \mathbf{W})$ . In addition, by (17b) and (17c), we have

$$\Delta \mathbf{w}^{T} \left( \mathbf{M}_{2}^{T} \otimes \mathbf{M}_{1} \right)^{T} \left( \mathbf{M}_{2}^{T} \otimes \mathbf{M}_{1} \right) \Delta \mathbf{w} = \Delta \mathbf{w}^{T} \left( \mathbf{M}_{2} \otimes \mathbf{M}_{1}^{T} \right) \left( \mathbf{M}_{2}^{T} \otimes \mathbf{M}_{1} \right) \Delta \mathbf{w}$$
$$= \Delta \mathbf{w}^{T} \left( \mathbf{M}_{2} \mathbf{M}_{2}^{T} \otimes \mathbf{M}_{1}^{T} \mathbf{M}_{1} \right) \Delta \mathbf{w}.$$

Finally, by exploiting the fact that  $\frac{\partial^2 \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}^2} = \mathbf{A} + \mathbf{A}^T$ , we obtain

$$\frac{\partial^2 \left\| \mathbf{M}_1 \Delta \mathbf{W} \mathbf{M}_2 \right\|_F^2}{\partial \Delta \mathbf{w}^2} = \mathbf{M}_2 \mathbf{M}_2^T \otimes \mathbf{M}_1^T \mathbf{M}_1 + \left( \mathbf{M}_2 \mathbf{M}_2^T \otimes \mathbf{M}_1^T \mathbf{M}_1 \right)^T$$
$$\stackrel{(a)}{=} \mathbf{M}_2 \mathbf{M}_2^T \otimes \mathbf{M}_1^T \mathbf{M}_1 + \left( \mathbf{M}_2 \mathbf{M}_2^T \right)^T \otimes \left( \mathbf{M}_1^T \mathbf{M}_1 \right)^T$$
$$= 2\mathbf{M}_2 \mathbf{M}_2^T \otimes \mathbf{M}_1^T \mathbf{M}_1,$$

where (a) follows from (17b). This completes the proof.

## B REFINED WEIGHT-UPDATE FORMULA

 We recall that the Hessian-based weight-update formula is given by (Frantar & Alistarh, 2022; Frantar et al., 2023)

$$\boldsymbol{\delta w} = -rac{w_q - \mathcal{Q}(w_q)}{[\mathbf{U}]_{q,q}} [\mathbf{U}]_{q,:}$$
 where  $\mathbf{U} = \mathrm{Chol}(\mathbf{H}^{-1})^T$ .

For the proposed attention-aware Hessians in Table 1, we have

$$\mathbf{U}_h = \mathbf{U}_{\mathrm{col},h} \otimes \mathbf{U}_{\mathrm{row},h},$$

where  $\mathbf{U}_{\text{col},h} = \text{Chol}(\mathbf{H}_{\text{col},h}^{-1})^T$  and  $\mathbf{U}_{\text{row},h} = \text{Chol}(\mathbf{H}_{\text{row},h}^{-1})^T$  (see Section 3.3). Therefore, the weight-update formula can be recast as

$$oldsymbol{\delta} oldsymbol{w}_h = -rac{w_q - \mathcal{Q}(w_q)}{[\mathbf{U}_{ ext{col},h} \otimes \mathbf{U}_{ ext{row},h}]_{q,q}} [\mathbf{U}_{ ext{col},h} \otimes \mathbf{U}_{ ext{row},h}]_{q,:}.$$



Figure 2: Illustration of the Hessian information when  $d_{row} = 2$  and  $d_{col} = 3$ 

For simplicity, suppose we quantize the first (0-th) row. When the weight  $[\mathbf{W}_h]_{0,j} (= [\mathbf{W}^{(0)}]_{h,j})$  in the *j*-th column is quantized, the weight-update of the *i*-th row is simplified as (see Fig. 2 for the ease of understanding)

$$\begin{split} [\boldsymbol{\delta}\mathbf{W}_{h}]_{i,:} &= -\frac{[\mathbf{W}_{h}]_{0,j} - \mathcal{Q}([\mathbf{W}_{h}]_{0,j})}{[\mathbf{U}_{\text{row},h}]_{0,0}[\mathbf{U}_{\text{col},h}]_{j,j}} [\mathbf{U}_{\text{row},h}]_{0,i}[\mathbf{U}_{\text{col},h}]_{j,:} \\ &= -\frac{[\mathbf{W}_{h}]_{0,j} - \mathcal{Q}([\mathbf{W}_{h}]_{0,j})}{[\mathbf{U}_{\text{col},h}]_{j,j}} \cdot \frac{[\mathbf{U}_{\text{row},h}]_{0,i}[\mathbf{U}_{\text{col},h}]_{j,:}}{[\mathbf{U}_{\text{row},h}]_{0,0}} \end{split}$$

Thus, after the quantization of all weights in the first row, the total amount of the weight-update for the *i*-th row can be expressed as

$$[\delta \mathbf{W}_{h, ext{total}}]_{i,:} = -\sum_{j=0}^{d_{ ext{col}}-1} rac{[\mathbf{W}_h]_{0,j} - \mathcal{Q}([\mathbf{W}_h]_{0,j})}{[\mathbf{U}_{ ext{col},h}]_{j,j}} \cdot rac{[\mathbf{U}_{ ext{row},h}]_{0,i}[\mathbf{U}_{ ext{col},h}]_{j,:}}{[\mathbf{U}_{ ext{row},h}]_{0,0}}$$

$$= -\frac{[\mathbf{U}_{\text{row},h}]_{0,i}}{[\mathbf{U}_{\text{row},h}]_{0,0}} \sum_{j=0}^{d_{\text{col}}-1} \frac{[\mathbf{W}_{h}]_{0,j} - \mathcal{Q}([\mathbf{W}_{h}]_{0,j})}{[\mathbf{U}_{\text{col},h}]_{j,j}} \cdot [\mathbf{U}_{\text{col},h}]_{j,:}.$$

Furthermore, by noting that (see line 8 in Algorithm 2)

$$[\mathbf{E}_{\text{GPTQ}}]_{h,j} = \frac{[\mathbf{W}_h]_{0,j} - \mathcal{Q}([\mathbf{W}_h]_{0,j})}{[\mathbf{U}_{\text{col},h}]_{j,j}}$$

we obtain

$$[\boldsymbol{\delta}\mathbf{W}_{h,\text{total}}]_{i,:} = -\frac{[\mathbf{U}_{\text{row},h}]_{0,i}}{[\mathbf{U}_{\text{row},h}]_{0,0}} \sum_{j=0}^{d_{\text{col}}-1} [\mathbf{E}_{\text{GPTQ}}]_{h,j} \cdot [\mathbf{U}_{\text{col},h}]_{j,:} = -\frac{[\mathbf{U}_{\text{row},h}]_{0,i}}{[\mathbf{U}_{\text{row},h}]_{0,0}} [\mathbf{E}_{\text{GPTQ}}]_{h,:} \mathbf{U}_{\text{col},h}.$$

As a result, the weight-update matrix to compensate for the quantization error of the first row is given by

$$[\boldsymbol{\delta}\mathbf{W}_{h,\text{total}}]_{0:,:} = -\frac{[\mathbf{U}_{\text{row},h}^{T}]_{0:,0}[\mathbf{E}_{\text{GPTQ}}]_{h,:}\mathbf{U}_{\text{col},h}}{[\mathbf{U}_{\text{row},h}]_{0,0}}.$$
(18)

By taking similar steps as above, we can easily generalize (18) for the *j*-th row as follows:  

$$[TT^T_{j} = 1 - TT_{j}]$$

$$[\boldsymbol{\delta}\mathbf{W}_{h,\text{total}}]_{j:,:} = -\frac{[\mathbf{U}_{\text{row},h}]_{j:,j}[\mathbf{E}_{\text{GPTQ}}]_{h,:}\mathbf{U}_{\text{col},h}}{[\mathbf{U}_{\text{row},h}]_{j,j}}.$$
(19)

#### С ADDITIONAL EXPERIMENTAL RESULTS

In this appendix, we provide experimental results omitted in the main text due to the page limitation.

#### C.1 COMPARISON WITH GPTQ

Table 7 and Table 8 summarize the INT3/INT4 quantization performances (perplexity) of the pro-posed BOA and the conventional GPTQ on various sizes of OPT, BLOOM, and LLaMA models. As evident from Table 3, Table 7, and Table 8, BOA uniformly outperforms GPTQ, and the performance gap is significant for low bit-width (i.e., INT2) and small-sized models suited for resource-limited devices (*e.g.*, mobile devices). 

Table 7: Quantization performance (PPL  $\downarrow$ ) of the proposed BoA and GPTQ on OPT.

			(a) Wi	kiText-2				
Precision	Method	125M	350M	1.3B	2.7B	6.7B	13B	30B
FP16	Baseline	27.65	22.00	14.63	12.47	10.86	10.13	9.56
	RTN	37.28	25.94	48.20	16.92	12.10	11.32	10.98
INT4	GPTQ	30.24	23.50	14.84	12.53	11.09	10.26	9.608
	BOA	28.93	22.90	14.72	12.44	10.88	10.16	9.571
	RTN	1.3e3	64.57	1.3e4	1.6e4	5.8e3	3.4e3	1.6e3
INT3	GPTQ	38.74	26.31	16.70	14.01	11.91	10.85	9.911
	BOA	33.68	24.69	15.93	13.43	11.53	10.58	9.826

(b) PTB

Precision	Method	125M	350M	1.3B	2.7B	6.7B	13B	30B
FP16	Baseline	38.99	31.08	20.29	17.97	15.77	14.52	14.04
INT4	RTN GPTQ BoA	53.88 44.31 <b>41.50</b>	36.79 33.41 <b>32.58</b>	75.37 21.23 <b>21.02</b>	32.41 18.70 <b>18.42</b>	18.86 16.09 <b>15.90</b>	16.41 14.69 <b>14.63</b>	15.44 14.18 14.18
INT3	RTN GPTQ <b>BoA</b>	1.4e3 57.62 <b>48.50</b>	87.21 39.35 <b>36.83</b>	1.5e4 24.77 <b>23.53</b>	1.4e4 21.53 <b>20.33</b>	5.3e3 17.56 <b>16.86</b>	2.2e3 15.68 <b>15.21</b>	1.5e3 14.56 <b>14.50</b>

(c) C4

Precision	Method	125M	350M	1.3B	2.7B	6.7B	13B	30B
FP16	Baseline	26.56	22.59	16.07	14.34	12.71	12.06	11.44
INT4	RTN GPTQ BoA	33.88 28.53 <b>27.56</b>	26.21 23.73 <b>23.20</b>	27.50 16.51 <b>16.39</b>	18.83 14.72 <b>14.61</b>	14.37 12.88 <b>12.84</b>	13.32 12.16 12.16	13.55 <b>11.50</b> 11.51
INT3	RTN GPTQ BoA	834.4 33.90 <b>31.12</b>	55.15 26.68 <b>25.39</b>	6.6e3 18.18 <b>17.74</b>	1.2e4 16.10 <b>15.83</b>	5.0e3 13.60 <b>13.34</b>	2.8e3 12.62 <b>12.52</b>	1.8e3 11.76 <b>11.75</b>

Table 8: Quantization performance (PPL  $\downarrow$ ) of BOA and GPTQ on BLOOM and LLaMA.

	(a) WikiText-2										
Dragision	Busician Mathed BLOOM										
Flecision	Method	560M	1.1B	1.7B	3B	7.1B	13B	30B			
FP16	Baseline	22.42	17.69	15.39	13.48	11.37	5.091	4.101			
	RTN	25.82	19.98	16.96	14.75	12.09	5.525	4.536			
INT4	GPTQ	23.44	18.54	15.90	13.90	11.63	5.262	4.285			
	BOA	23.28	18.32	15.81	13.84	11.58	5.243	4.262			
	RTN	56.74	49.85	63.37	39.07	17.35	11.78	14.87			
INT3	GPTQ	26.63	20.80	17.71	15.39	12.42	5.721	4.848			
	BOA	25.90	20.28	17.12	14.91	12.19	5.676	4.725			

### (b) PTB

Precision	Method		LL	LLaMA				
TICCISION	Wiethou	560M	1.1B	1.7B	3B	7.1B	13B	30B
FP16	Baseline	43.69	57.96	30.00	25.34	20.83	9.081	8.159
	RTN	50.96	66.79	33.52	27.65	22.40	9.775	8.653
INT4	GPTQ	45.33	61.94	31.37	26.39	21.40	9.306	8.344
	BOA	44.92	61.40	30.67	26.23	21.34	9.255	8.304
	RTN	124.8	184.0	105.5	66.24	34.94	28.94	28.79
INT3	GPTQ	52.39	70.68	35.06	28.99	23.46	9.928	8.925
	BOA	50.71	67.77	33.92	28.67	22.86	9.857	8.737

# (c) C4

Precision	Method		BLOOM						
Treeision	Wiethou	560M	1.1B	1.7B	3B	7.1B		13B	30B
FP16	Baseline	26.60	22.05	19.49	17.49	15.20		6.798	6.131
	RTN	29.80	24.42	21.24	18.75	16.05		7.232	6.537
INT4	GPTQ	27.39	22.69	20.03	17.89	15.44		6.973	6.294
	BOA	27.23	22.54	19.90	17.82	15.42		6.958	6.267
	RTN	66.99	60.41	113.6	79.84	22.54		14.46	30.04
INT3	GPTQ	29.89	24.48	21.44	19.07	16.24		7.504	6.840
	BOA	29.39	24.17	21.02	18.74	16.09		7.454	6.718

### 918 C.2 INTEGRATION WITH EQUIVALENT TRANSFORM-BASED METHODS

In this appendix, we verify that the performance of the proposed BOA can be enhanced by combining BOA with existing ET-based methods. Among various algorithms, we use SmoothQuant (Xiao et al., 2023) and Z-FOLD (Jeon et al., 2023b) in our integration because they efficiently find out an equivalent transform without time-consuming gradient-based optimization. We note that while SmoothQuant has been proposed in the context of weight-activation quantization, the smoothing factor  $(s_j = \max(|\mathbf{X}_j|)^{\alpha} / \max(|\mathbf{W}_j)^{1-\alpha})$  used for the equivalent transformation can also be used for weight-only quantization by setting  $\alpha = 0$ .

Tables 5 and 9 summarize the integration results. Overall, the performances of BOA and GPTQ
indeed improve when combined with ET-based methods. We emphasize that the performance gap
between the proposed BOA and GPTQ still remains significant, especially for INT2 quantization. A
similar behavior can be observed in the zero-shot results (see Table 10); the performance is boosted
by applying ET-based methods, and BOA outperforms GPTQ regardless of the ET-based method.

Table 9: INT3 performance (PPL  $\downarrow$ ) of BOA integrated with existing ET-based methods.

Equivalent Transformation	Dataset	Method	125M	350M	Mod 1.3B	el Size (C 2.7B	0PT) 6.7B	13B	30B
	Wiki2	GPTQ BoA	39.56 <b>34.58</b>	N/A N/A	16.32 <b>15.83</b>	13.55 <b>13.33</b>	11.90 <b>11.58</b>	10.68 <b>10.38</b>	9.857 <b>9.846</b>
SmoothQuant	РТВ	GPTQ BoA	58.00 <b>51.44</b>	N/A N/A	24.00 <b>22.78</b>	20.36 <b>19.83</b>	17.18 <b>16.74</b>	15.45 <b>15.18</b>	14.46 <b>14.40</b>
	C4	GPTQ BoA	34.98 <b>31.52</b>	N/A N/A	17.78 <b>17.43</b>	15.68 <b>15.48</b>	13.50 13.35	12.55 <b>12.48</b>	11.74 <b>11.73</b>
	Wiki2	GPTQ BoA	39.59 <b>33.31</b>	25.97 <b>24.22</b>	16.10 <b>15.91</b>	13.54 <b>13.33</b>	11.65 <b>11.28</b>	10.64 <b>10.53</b>	9.887 <b>9.814</b>
Z-Fold	РТВ	GPTQ BoA	53.08 <b>46.59</b>	39.23 <b>36.80</b>	22.73 <b>22.29</b>	20.18 <b>19.54</b>	16.64 <b>16.44</b>	15.22 <b>15.16</b>	14.57 <b>14.53</b>
	C4	GPTQ BoA	33.67 <b>30.00</b>	26.45 <b>25.04</b>	17.33 <b>17.13</b>	15.50 <b>15.32</b>	13.28 13.20	12.46 <b>12.41</b>	11.73 <b>11.71</b>

\* SmoothQuant does not support OPT-350M where the post-LayerNorm architecture has been used.

Table 10: INT2 zero-shot performance (accuracy ↑) of BOA integrated with ET-based methods.

Equivalent	Model Size	Mathod			Tasks		Avorea
Transformation	(OPT)	Method	ARC-c	ARC-e	HellaSwag	MMLU	Average
	1 2 D	GPTQ	23.38	40.15	37.47	23.00	31.00
	1.5D	BOA	22.18	42.00	37.48	23.10	31.19
	2.70	GPTQ	25.85	42.72	42.46	23.10	33.53
	2./B	BOA	27.73	44.70	44.24	22.95	34.91
SmoothQuant	6 7D	GPTQ	25.68	46.25	45.24	23.40	35.14
ShioothQuant	0./B	BOA	27.56	48.15	46.20	23.90	36.45
	12D	GPTQ	29.52	52.53	56.63	24.90	40.90
	13B	BOA	31.57	53.66	58.69	25.35	42.32
	200	GPTQ	29.18	54.46	60.04	24.54	42.06
	30B	BOA	31.57	55.93	62.18	25.54	43.81
	1.2D	GPTQ	23.89	42.05	40.45	23.07	32.37
	1.3B	BOA	24.74	43.98	40.31	23.45	33.12
	2 7B	GPTQ	25.00	41.46	43.13	23.16	33.19
	2.7 <b>D</b>	BOA	26.37	43.06	45.18	23.50	34.53
7 FOLD	6.7D	GPTQ	30.46	48.78	52.46	25.64	39.34
Z-POLD	0./B	BOA	28.58	49.75	55.45	26.67	40.11
	12D	GPTQ	28.58	48.78	55.38	24.75	39.37
	13B	BOA	28.84	49.87	58.32	24.60	40.41
	200	GPTQ	31.83	53.70	61.34	24.96	42.96
	20B	BOA	30.12	57.53	63.63	24.85	44.03

#### C.3 COMPARISON WITH PRIOR ARTS

We compare the proposed BOA with OmniQuant (Shao et al., 2023) and AffineQuant (Ma et al., 2024), recently proposed algorithms that learn an attention-aware equivalent transform via back-propagation (see Tables 6 and 11 for PPL results and see Table 12 for zero-shot results).

As evident, BOA itself outperforms existing algorithms in almost all cases, even though BOA does not rely on time-consuming gradient-based optimization. Furthermore, when combined with SmoothQuant or Z-FOLD, the performance gap between BOA and OmniQuant/AffineQuant is sig-nificant, which demonstrates the efficacy of BOA. We note that OmniQuant and AffineQuant sometimes diverge or collapse (*i.e.*, PPL is larger than  $10^3$ ) for INT2 quantization. In fact, to supplement the INT2 quantization performance, group-wise quantization parameters have been additionally used in OmniQuant and AffineQuant, but group-wise parameters result in additional memory costs and processing time for the inference (Shen et al., 2023). 

Table 11: INT3 performance (PPL  $\downarrow$ ) of BOA and existing attention-aware approaches.

Dotocot	Mathod		Model Size (OPT)							
Dataset	Wiethou	125M	1.3B	2.7B	6.7B	13B	30B			
	OmniQuant	41.59	18.23	15.11	12.86	12.49	11.26			
	AffineQuant	37.75	17.12	14.32	12.42	11.93	10.72			
Wiki2	воА	33.68	15.93	13.43	11.53	10.58	9.826			
	BOA + SmoothQuant	34.58	15.83	13.33	11.58	10.38	9.846			
	BOA + Z-FOLD	33.31	15.91	13.33	11.28	10.53	9.814			
	OmniQuant	59.51	26.08	22.68	18.31	17.76	16.02			
	AffineQuant	53.02	24.47	21.18	17.27	17.27	15.31			
PTB	воА	48.50	23.53	20.33	16.86	15.21	14.50			
	BOA + SmoothQuant	51.44	22.78	19.83	16.74	15.18	14.40			
	BOA + Z-FOLD	46.59	22.29	19.54	16.44	15.16	14.53			
	OmniQuant	35.73	19.10	16.80	14.40	13.48	12.44			
	AffineQuant	33.37	18.56	16.15	13.91	13.31	12.14			
C4	воА	31.12	17.74	15.83	13.34	12.52	11.75			
	BOA + SmoothQuant	31.52	17.43	15.48	13.35	12.48	11.73			
	BOA + Z-FOLD	30.00	17.13	15.32	13.20	12.41	11.71			

Table 12: Zero-shot task performance (accuracy↑) of BOA and existing attention-aware methods

1004	Model Size			r	Fasks		
1005	(OPT)	Method	ARC-c	ARC-e	HellaSwag	MMLU	Average
1006		OmniQuant	NaN	NaN	NaN	NaN	NaN
1007	1 2 P	AffineQuant	NaN	NaN	NaN	NaN	NaN
1008	1.5D	BOA + SmoothQuant	22.18	42.00	37.48	23.10	31.19
1009		BOA + Z-FOLD	24.74	43.98	40.31	23.45	33.12
1010		OmniQuant	NaN	NaN	NaN	NaN	NaN
1011	2 7B	AffineQuant	25.09	40.66	39.35	22.90	32.00
1012	2.75	BOA + SmoothQuant	27.73	44.70	44.24	22.95	34.91
1013		BOA + Z-FOLD	26.37	43.06	45.18	23.50	34.53
1014		OmniQuant	23.55	28.87	25.60	22.95	25.24
1015	6 7B	AffineQuant	26.62	48.23	47.18	23.57	36.40
1016	0.72	BOA + SmoothQuant	27.56	48.15	46.20	23.90	36.45
1017		BOA + Z-FOLD	28.58	49.75	55.45	26.67	40.11
1012		OmniQuant	26.11	26.35	25.54	22.95	25.24
1010	13B	AffineQuant	24.15	43.48	45.62	22.89	34.04
1019	102	BOA + SmoothQuant	31.57	53.66	58.69	25.35	42.32
1020		BOA + Z-FOLD	28.84	49.87	58.32	24.60	40.41
1021		OmniQuant	26.11	26.73	25.83	22.95	25.41
1022	200	AffineQuant	25.26	26.47	25.65	22.95	25.08
1023	30B	BOA + SmoothQuant	31.57	55.93	62.18	25.54	43.81
1024		BOA + Z-FOLD	30.12	57.53	63.63	24.85	44.03
1025	* 'NaN' mea	ns that loss diverges in the	auantizatio	n process			

'NaN' means that loss diverges in the quantization process.

# 1026 C.4 COMPARISON OF TIME AND MEMORY COSTS

In this appendix, we compare the processing time and memory costs of BOA and conventional algorithms.

Table 13: Time and memory costs of the proposed BOA and existing methods

(a) INT2 quantization processing time

Mathod	Reconstruction	Model Size (OPT)							
Method	Target	125M	1.3B	2.7B	6.7B	13B	30B		
GPTQ	Layer output	0.752 min	6.284 min	0.214 hr	0.603 hr	1.293 hr	3.689 hr		
OmniQuant	Transformer block output	16.20 min	61.20 min	1.627 hr	2.933 hr	5.309 hr	11.57 hr		
AffineQuant	Transformer block output	28.33 min	154.2 min	4.597 hr	9.854 hr	18.41 hr	44.25 hr		
BoA	Attention output	5.099 min	31.64 min	1.101 hr	2.830 hr	4.964 hr	10.55 hr		
		(b) Memo	orv cost (GB	)					

Mathad	Reconstruction			Model Si	ze (OPT)		
Method	Target	125M	1.3B	2.7B	6.7B	13B	
GPTQ	Layer output	1.391	3.970	4.863	9.011	12.57	2
OmniQuant	Transformer block output	1.941	5.869	7.095	11.68	16.15	2
AffineQuant	Transformer block output	3.473	9.963	12.25	20.08	26.78	4
BOA	Attention output	1.676	5.471	6.837	12.26	19.06	3

**Discussion on BOA** The main reason why the processing time of the proposed BOA increases with the size of LLMs is that the embedding dimension of attention heads increases with the model size. Specifically, to compensate for the error incurred by the quantization of certain rows, BOA needs to sequentially quantize sub-weight matrix  $d_h$  times where  $d_h$  is the head dimension (see Fig. 1(b)). Because  $d_h$  increases with model size, the number of sequential quantizations also increases, which leads to long processing time.

1055

1030

1031 1032

1033 1034 1035

1039 1040

1041

1043

1056 **Comparison with OmniOuant/AffineOuant** The processing time of the proposed BOA is shorter 1057 than those required by existing attention-aware algorithms that rely on gradient-based optimization, 1058 yet achieving significantly better performance (see Tables 6 and 11). We note that OmniQuant does 1059 not take too much time for quantizing large LLMs even though it performs gradient-based optimization. This is because OmniQuant reduces the number of learnable parameters greatly to accelerate 1061 gradient-based optimization. Specifically, instead of learning a weight-rounding policy which requires to learn a large number of parameters, OmniQuant learns only a small number of quantiza-1062 tion parameters and some parameters related to the equivalent transform (Shao et al., 2023). While 1063 this strategy accelerates the quantization process greatly, OmniQuant suffers from unstable quanti-1064 zation process or collapses for low-bit quantization (see Tables 6 and 11). AffineQuant improves OmniQuant by introducing additional learnable parameters (Ma et al., 2024), but such additional 1066 parameters result in huge processing time (4 times longer processing time; see Table 13(a)), which 1067 demonstrates the inefficiency of gradient-based optimization over the proposed method. 1068

1069

**Comparison with GPTQ** The proposed BOA requires longer processing time and larger memory 1070 than those required by GPTQ. This is because GPTQ quantizes all the rows of the weight matrix 1071 simultaneously using only layer input. In contrast, BOA sequentially quantizes sub-weight matri-1072 ces using outputs of other layers as well as layer input (see Fig. 1(b)) to consider the inter-layer dependencies within the attention module, which eventually leads to the better quantization per-1074 formance than GPTQ. Clearly, there is a trade-off between quantization speed / memory cost and 1075 accuracy. In real situations, when one needs to preserve the performance of the original model by considering inter-layer dependencies within the attention module, the proposed BoA would be an intriguing solution. Even when the memory resource is limited, BOA can be used with some relax-1077 ation. Specifically, we note that the large memory cost of BOA for hyper-scale LLMs (e.g., 13B and 1078 30B) is attributable to the row-wise Hessian for the value projection  $(\mathbf{X}\mathbf{A}_h^T\mathbf{A}_h\mathbf{X}^T$  in (12)) whose 1079 shape is  $H \times d \times d$  (H is the number of attention heads and d is the embedding dimension). In

memory-limited cases, we can mitigate the memory cost of BOA by considering inter-layer depen-1081 dencies only for query and key projections and applying the standard Hessian ( $\mathbf{X}\mathbf{X}^{T}$  in (4)) for the 1082 value projection. Indeed, when considering only query and key projections, BoA requires almost 1083 same amount of memory as GPTQ, yet still exhibiting better performance (see Table 14).

Table 14: Performance and memory costs of the proposed BOA with and without considering inter-1085 layer dependencies for the value projection 1086

> Consideration of Model Size (OPT) Method Inter-layer Dependencies 125M 1.3B 2.7B 6.7B 13B 30B 33.94 24.86 20.08 14.45 GPTQ 178.6 64.11 OmniQuant query, key, value, out, fcs NaN NaN NaN 3.0e4 2.0e5 2.1e5 107.0 34.45 31.50 AffineQuant query, key, value, out, fcs NaN 25.11 3.3e5 BOA 130.8 52.84 27.81 23.95 19.98 14.11 query, key 48.92 26.57 23.03 19.22 BOA query, key, value 118.1 13.84

(a) INT2 p	erformance	(PPL J)	on C <sup>2</sup>
------------	------------	---------	-------------------

### (b) Memory cost (GB)

Method	Consideration of Inter-layer Dependencies	125M	1.3B	Model Si 2.7B	ze (OPT) 6.7B	13B	30B
GPTQ OmniQuant AffineQuant	query, key, value, out, fcs query, key, value, out, fcs	1.391 1.941 3.473	3.970 5.869 9.963	4.863 7.095 12.25	9.011 11.68 20.08	12.57 16.15 26.78	21.25 26.94 42.21
BOA BOA	query, key query, key, value	1.391 1.676	3.971 5.471	4.864 6.837	9.015 12.26	12.57 19.06	21.25 39.74

The additional memory required for the query and key projections, i.e., memory needed to save  $\mathbf{K}_{h}^{T}\mathbf{K}_{h}$  and  $\mathbf{Q}_{h}^{T}\mathbf{Q}_{h}$  (see (15) and (16)), is negligible (e.g., 0.003 GB for 30B).

1080

1084

1087 1088 1089

1090

1091

1092

1093

1094

1095 1096

# 1134 C.5 RESULTS FOR DIFFERENT CALIBRATION DATASETS

When constructing a calibration dataset, we randomly sample 128 sequences from the C4 dataset (see Section 4.1). By changing the seed for the sampling, we can obtain different calibration datasets, which leads to different quantization results.<sup>7</sup> In this appendix, we report the corresponding results and overall statistics. Due to the limited computational resources, we conducted this experiment only for our main comparison (*i.e.*, the performances of the proposed BoA and the conventional GPTQ).

1142 1143

Table 15: Performance (perplexity  $\downarrow$ ) of the proposed BOA and GPTQ for different seeds.

				(a) INT2 (	Quantization			
Dataset	Seed	Method	125M	1.3B	2.7B	6.7B	13B	30
	0	GPTQ BoA	232.8 <b>141.6</b>	66.76 <b>48.71</b>	37.44 <b>26.20</b>	24.74 <b>22.71</b>	18.97 <b>18.76</b>	13. 12.
Wiki2	10	GPTQ BoA	276.2 147.4	66.30 <b>47.23</b>	36.74 <b>25.95</b>	24.64 23.11	20.05 18.52	13. 12.
	20	GPTQ BoA	243.9 <b>139.4</b>	65.09 <b>45.11</b>	36.33 <b>27.00</b>	24.94 <b>24.06</b>	19.78 <b>17.82</b>	13. 12.
	50	GPTQ BoA	269.9 <b>160.7</b>	64.88 <b>44.13</b>	33.84 <b>26.75</b>	24.54 <b>23.09</b>	19.47 <b>18.70</b>	13. 12.
	100	GPTQ BoA	228.3 <b>147.8</b>	71.51 <b>47.43</b>	36.72 <b>26.85</b>	25.55 <b>23.61</b>	19.39 <b>18.49</b>	13. 12.
	Mean ±Stdev	GPTQ BoA	250.2±22 <b>147.4±8.3</b>	66.91±2.7 <b>46.52</b> ±1.9	36.21±1.4 26.55±0.45	24.88±0.40 23.32±0.53	19.53±0.41 <b>18.46±0.38</b>	13.24: <b>12.19</b>
	0	GPTQ BoA	384.8 <b>199.2</b>	112.0 <b>78.73</b>	64.59 <b>40.76</b>	42.36 <b>33.77</b>	26.95 <b>25.34</b>	20. 18.
PTB	10	GPTQ BoA	324.1 <b>185.4</b>	112.7 <b>76.02</b>	62.42 <b>39.73</b>	38.91 <b>33.79</b>	27.92 23.55	19. <b>18.</b>
	20	GPTQ BoA	350.4 <b>188.7</b>	111.6 <b>79.31</b>	62.64 <b>41.69</b>	39.84 <b>34.89</b>	27.80 <b>24.59</b>	20. <b>18</b> .
	50	GPTQ BoA	433.8 <b>206.8</b>	122.1 <b>87.29</b>	59.46 <b>41.91</b>	43.05 <b>36.36</b>	27.64 <b>24.56</b>	20 18
	100	GPTQ BoA	479.2 <b>164.8</b>	125.8 77.32	59.49 <b>41.67</b>	38.26 <b>35.19</b>	27.56 <b>24.56</b>	20. 17.
	Mean ±Stdev	GPTQ BoA	394.5±63 <b>189.0</b> ± <b>16</b>	116.8±6.6 <b>79.73</b> ± <b>4.4</b>	61.72±2.2 <b>41.15±0.91</b>	40.48±2.1 <b>34.80±1.1</b>	27.57±0.38 24.52±0.64	20.21 18.18
	0	GPTQ BoA	178.6 <b>118.1</b>	64.11 <b>48.92</b>	33.94 <b>26.57</b>	24.86 23.03	20.08 <b>19.22</b>	14. <b>13.</b>
	10	GPTQ BoA	189.7 <b>115.8</b>	64.19 <b>49.76</b>	33.27 <b>27.00</b>	24.40 <b>24.04</b>	20.40 <b>18.58</b>	14. <b>13</b> .
C4	20	GPTQ BoA	163.1 <b>111.1</b>	64.19 <b>48.81</b>	32.83 <b>26.96</b>	24.66 23.31	20.37 <b>19.15</b>	14. 13.
C4	50	GPTQ BoA	190.8 <b>119.2</b>	64.99 <b>46.68</b>	32.69 <b>27.11</b>	24.55 23.73	20.13 <b>18.96</b>	14. <b>13.</b>
	100	GPTQ BoA	168.5 <b>116.7</b>	67.54 <b>48.49</b>	33.65 27.17	25.25 <b>24.09</b>	20.23 <b>19.63</b>	14 <b>13</b>
	Mean ±Stdev	GPTQ BoA	178.2±12 116.2±3.1	65.00±1.5 <b>48.53</b> ± <b>1.1</b>	33.28±0.53 26.96±0.24	24.75±0.33 23.64±0.46	20.24±0.14 19.11±0.38	14.46± 13.83±

1181 1182 1183

1180

1184

1185

1186 1187

<sup>7</sup>Tables 3 to 11 present the results for seed 0.

				(b) INT3 Q	uantization			
Dataset	Seed	Method	125M	1.3B	2.7B	6.7B	13B	30B
	0	GPTQ	38.74	16.70	14.01	11.91	10.85	9.911
		BOA	33.68	15.93	13.43	11.53	10.58	9.826
	10	GPTQ	40.72	16.99	13.97	11.89	10.84	9.881
		BUA	34.30	10.15	13.43	11.43	10.59	9.750
	20	GPTQ BoA	38.72 34.24	17.17 16.02	13.89 13.32	11.98 <b>11.66</b>	10.94 10.51	9.941 <b>9.791</b>
Wiki2		GPTQ	37.72	16.93	14.08	11.86	10.92	9.769
	50	BOA	34.73	16.19	13.44	11.44	10.54	9.768
	100	GPTQ	38.37	16.97	14.27	11.89	10.88	9.840
		BOA	34.51	16.14	13.72	11.49	10.52	9.847
	Mean	GPTQ	$38.85 \pm 1.1$	$16.95 \pm 0.17$	$14.04 \pm 0.14$	$11.91 \pm 0.044$	$10.89 \pm 0.045$	$9.868 \pm 0.0$
	±Stdev	воа	34.29±0.39	16.09±0.11	13.47±0.15	11.51±0.094	$10.55 \pm 0.036$	9.798±0.
	0	GPTQ Boa	57.62 48 50	24.77 23.53	21.53 20 33	17.56 16 86	15.68 15 21	14.56 14 50
		CDTO	<b>40.50</b>	25.35	20.55	17.10	15.21	14.50
	10	BOA	<b>48.37</b>	23.19 23.35	21.60 20.27	17.19 16.71	15.00 15.13	14.00 14.50
		GPTO	56.22	25.47	21.42	17.24	15.67	14 60
	20	BOA	47.01	23.31	20.22	16.68	15.16	14.54
PTB	50	GPTQ	56.68	24.94	21.44	17.36	15.74	14.52
	50	BOA	51.05	23.69	20.19	16.76	15.22	14.46
	100	GPTQ	52.15	25.13	21.32	17.44	15.79	14.53
		BOA	48.66	23.50	19.93	16.77	15.21	14.46
	Mean ⊥Stdov	GPTQ	$55.65 \pm 2.1$	$25.10 \pm 0.26$	$21.46 \pm 0.11$	$17.36 \pm 0.15$	$15.71 \pm 0.057$	$14.57 \pm 0.$
	Tantes	DUA	40.72 ± 1.5	23.48±0.13	20.19±0.10	10.75±0.008	13.19±0.038	14.49±0.
	0	BOA	33.90 <b>31.12</b>	18.18 17.74	16.10 15.83	13.60 13.34	12.62 12.52	11.76
		GPTO	34.16	18 19	16.07	13.60	12.63	11.76
	10	BOA	31.72	17.72	15.70	13.34	12.52	11.74
		GPTQ	34.07	18.19	16.07	13.58	12.62	11.75
C1	20	BOA	31.29	17.72	15.73	13.33	12.53	11.74
C4	50	GPTQ	33.68	18.16	16.06	13.60	12.67	11.76
		BOA	31.11	17.73	15.73	13.37	12.54	11.74
	100	GPTQ	33.80	18.20	16.12	13.61	12.66	11.76
		BOA	31.38	17.75	15./0	13.30	12.55	11.75
	Mean +Stdev	GPTQ Boa	$33.92 \pm 0.20$ $31.32 \pm 0.25$	18.18±0.015 17.73+0.013	16.08±0.026 15.74+0.053	$13.60 \pm 0.014$ $13.35 \pm 0.016$	$12.64 \pm 0.023$ 12.53+0.012	11.76±0.0
	Touch	DUA	01.02 ± 0.20	1	10.74 ± 0.000	10.00 ± 0.010	12.00 ± 0.012	11.75±0.0

# Under review as a conference paper at ICLR 2025

### 1242 D PSEUDOCODE FOR GPTQ 1243

In this appendix, we provide the pseudocode of the conventional GPTQ (Frantar et al., 2023), which is omitted in the main manuscript due to the page limitation.

1247 Al	gorithm 2 GPTQ
1248 In	<b>put</b> : weights W, Hessian information $U_{col}$ , pre-determined step size S, and blocksize B
1249	the def GPTQ(W, $U_{col}$ , S, $B = 128$ )
1250	: Initialize quantized output: $\mathbf{Q} \leftarrow 0_{d_{\text{rev}} \times d_{\text{rel}}}$
1251 3	: Initialize total quantization errors: $\mathbf{E}_{\text{total}} \leftarrow 0_{d_{\text{error}} \times d_{\text{rel}}}$
1252 4	Initialize block quantization errors: $\mathbf{E}_{block} \leftarrow 0_{d_{row} \times B}$
1253	: for $i = 0, B, 2B, \dots$ do
1254 6	for $j = i, \dots, i + B - 1$ do
1255	: Ouantize the <i>j</i> -th column: $\mathbf{Q}_{i} \leftarrow \text{quant}(\mathbf{W}_{i}, \mathbf{S})$
8	Estimate quantization error: $[\mathbf{E}_{block}]_{i,j=i} \leftarrow (\mathbf{W}_{i,j} - \mathbf{Q}_{i,j})/[\mathbf{U}_{col}]_{i,j}$
257 950 9	: Update weights in block: $\mathbf{W}_{:,i;i+B} \leftarrow \mathbf{W}_{:,i;i+B} - [\mathbf{E}_{block}]_{:,i-i} \cdot [\mathbf{U}_{col}]_{i,i:(i+B)}$
10	: end for
9 11	Update all remaining weights: $\mathbf{W}_{i+B} \leftarrow \mathbf{W}_{i+B} - \mathbf{E}_{block} \cdot [\mathbf{U}_{col}]_{i:(i+B)}$
12	: Save block quantization errors: $[\mathbf{E}_{total}] \cdot i \cdot i + B \leftarrow \mathbf{E}_{block}$
13	end for
	<b>utput</b> : quantized weights $\mathbf{Q}$ , quantization error $\mathbf{E}_{total}$
30 <u>-</u>	
04 CE	
60	
:0	
0	
9	
U 	
0	
2	
Л	
+	
7	
2	
, )	
- )	
)	
-	
r L	
, )	
1	
2	
2	
<u>л</u>	
10	