Achieving binary weight and activation for LLMs using Post-Training Quantization

Anonymous ACL submission

Abstract

Quantizing large language models (LLMs) to 002 1-bit precision significantly reduces computational costs, but existing quantization techniques suffer from noticeable performance degradation when using weight and activation precisions below 4 bits (W4A4). In this paper, 007 we propose a post-training quantization framework specifically designed for LLMs, utilizing a novel $W(1+1)A(1\times 4)$ configuration, where weights are quantized to 1 bit and activations are quantized to 1 bit with a 4-fold increase in the number of channels. For weight quantization, we propose utilizing Hessian-aware 013 fine-grained grouping along with an EM-based quantization scheme. For activation quantization, we decompose INT4-quantized activations into a 4 × INT1 format equivalently and si-017 multaneously smooth the scaling factors based on quantization errors, which further reduces the quantization errors in activations. This ef-021 fectively reduces activation quantization errors and enables the use of INT1 multiplication to accelerate overall inference computations. Our method surpasses state-of-the-art (SOTA) LLM quantization baselines on W2A4 across multiple tasks, pushing the boundaries of existing LLM quantization methods towards fully binarized models.

1 Introduction

037

041

The enormous computational and memory overhead of large language models (LLMs) limits their widespread adoption. Model quantization methods (Nagel et al., 2021; Huang et al., 2024b) is a major approach to alleviate this challenges. Recent studies (Zhao et al., 2024) on post-training quantization (PTQ) techniques can attain nearly lossless quantization under the W4A4 configuration, facilitating accelerated computation through the utilization of the INT4 format. However, lower-bit quantization of LLMs remains an unsolved problem when the bit number of both weight and activation is less



Figure 1: The perplexity (\downarrow) of GPTQ, QuaRot, Atom, and our method on the LLAMA1-7B model under various quantization bit-width settings reveals that GPTQ, QuaRot, and Atom exhibit significant performance degradation when weights are quantized to 1 bit, whereas our method still demonstrates language generation capabilities close to those of the original model.

than 4. As shown in Figure 1, existing methods (Zhao et al., 2024; Ashkboos et al., 2024) still cannot resolve the performance collapse that occurs when weights are represented using binarization.

043

045

047

048

051

054

059

060

061

062

Most PTQ LLM quantization work adopts the Round To Nearest (RTN) method (Nagel et al., 2021) as the implementation scheme in practical operations, majorly focused on how to preprocess data for RTN quantization, i.e., how to process the dataset before RTN quantization to make it easier to quantize. For example, LLM.int8 (Dettmers et al., 2022) and Atom (Zhao et al., 2024) handle outliers with high precision, effectively increasing the numerator of the scaling factor for some key values, thereby reducing the difficulty of quantization. An other line of Works, including SmoothQuant (Xiao et al., 2023), QuaRot (Ashkboos et al., 2024), and DuQuant (Lin et al., 2024) smooth the distribution of data to be quantized via matrix rotation, to reduce the denominator of the scaling factor and also lowering the quantization difficulty. Although

063

095

- 100
- 101
- 103
- 104

105

106 107 108

111

112

109 110

these works can achieve excellent 4-bit and 8-bit quantization, they struggle to further reduce the quantization bit-width. This problem largely stems from the fact that the simple RTN quantization method itself is difficult to support the quantization quality required for very low-bit quantization.

An other line of works use vector Ouantization, such as VPTQ (Liu et al., 2024a) and GPTVQ (van Baalen et al., 2024). This method clusters the data to be quantized into several classes and maps the data in the same class to the value with the smallest total distance. The quantization accuracy of this method is often higher than that of the RTN method at the same level. However, since the quantized values obtained by this method are not equally spaced, inference cannot be accelerated using efficient low-bit arithmetic with the quantized values. Instead, dequantization must be performed before inference, which negates all the speed-up benefits of the quantization method.

To address these issues, we propose a quantization framework, that can simultaneously achieve small model size, fast computation, and highquality inferece outputs. It consists of three parts: 1-bit weight quantization with another bit for finegrain grouping, and 1×4 bits activation quantization using 4 times 1-bit channel to represent 4bit quantization. By this way, we are able to compute the inner loop vector product in pure binary operations drastically boosting the inference speed as well as reduce the model size. The main contributions of this work are as follows:

- We propose a novel $W(1+1)A(1\times 4)$ posttraining quantization framework, where weights are quantized to 1 bit with additional bit for fine-grain grouping and activations are quantized to 1 bit with a 4-fold increase in the number of channels.
- We present an EM-based algorithm for searching (1+1) bit weight quantization, noticeably improved the performance compared with the usual RTN approach wildly used in other PTQ methods with fast inference speed and smaller model size.
- We compare our proposed method with the current state-of-the-art quantization methods. Experimental results show that our method significantly outperforms existing methods, especially at very low bit-widths. On the Wikitext2 benchmark, our method achieves

perplexities of 8.58 and 8.89 on LLaMA-7B and LLaMA2-7B, respectively, using only the W1A4 quantization setting. This significantly surpasses the existing state-of-the-art methods and is comparable to the performance of 5.68 and 5.47 achieved by the original fullprecision models.

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

154

155

156

157

158

159

160

161

162

Related Works 2

In recent years, considerable research efforts have been devoted to the low-bit quantization of LLMs. Current works have achieved near-lossless quantization at W4A4 through methods such as mixedprecision quantization (Zhao et al., 2024) and outlier smoothing using rotation factors (Ashkboos et al., 2024; Liu et al., 2024b). However, these methods encounter difficulties when attempting to push towards even lower-bit quantization.

Consequently, some studies have begun to focus on binarization methods (Courbariaux and Bengio, 2016; Qin et al., 2022) for LLMs. Most of these works are based on Quantization-Aware Training (QAT) theory. BitNet b1.58 (Wang et al., 2023), BitNet a4.8 (Wang et al., 2024), OneBit (Xu et al., 2024), and FBI-LLM (Ma et al., 2024) have designed 1-bit Transformer architectures specifically for LLMs, replacing all linear layers in the original LLMs with specific quantized linear layers. These methods can reduce the quantization bit-width of weights to 1 bit and demonstrate significant advantages over baseline models. Nevertheless, these QAT methods requires substantial computational resources, making it impractical for the quantization of LLMs.

Some studies have also explored the possibility of binarizing LLMs within the PTQ framework. BiLLM (Huang et al., 2024a) and STBLLM (Dong et al., 2024), inspired by the distribution characteristics of weight values and the Hessian matrix, adopt binary residual approximations for significant weights and optimal partitioning for group binarization of non-significant weights, pushing the quantization boundary of LLM weights down to 1 bit. However, most of these works exhibit suboptimal performance and necessitate an additional bit to store fine-grained grouping information. Moreover, these works neglect the quantization of activation values, which obstructs the computational acceleration of the quantized model and ultimately results in relatively poor performance.

Both ABQ-LLM (Zeng et al., 2024) and QBB



Figure 2: binarized weight and activation attention

(Bulat et al., 2024) propose the idea of decomposing a high-bit matrix into a set of binary matrices for accelerated computation. These methods hold promise for achieving INT1 computational acceleration but are constrained by the limitations of quantization bit-width and performance, preventing further improvements.

Based on current research on LLM binarization, we believe that weight binarization of LLMs within the PTQ framework is feasible. Compared to the current state-of-the-art (SOTA) W4A4 quantization efforts, $W(1+1)A(1\times4)$ quantized models can further reduce computational and memory overhead. Therefore, we propose our method, which further pushes the boundaries of LLM quantization and enriches the research on LLM binarization.

3 Method

163

164

165

166

167

168

169

171

172

173

174

175

176

177

178

179

180

181

185

186

187

189

191

192

193

In this section, we introduce our Binarized Weight and Activation (BWA) approach, a novel paradigm designed to simultaneously accelerate inference speed and reduce memory footprint in Large Language Models (LLMs). We begin by outlining the architectural framework of the proposed BWA attention module. Subsequently, we describe an EMbased method for computing binarized weights that are compatible with the proposed two-level grouping structure. Finally, we discuss the strategies employed in this study and compare them to other relevant methods.

3.1 Binarized weight and activation (BWA) attention

194 Overall structure of BWA attention Our BWA195 framework is a quantization-aware modification of

the standard attention module of a LLAMA-like model.

Following the dataflow of the activations in Figure 2.(a), we introduce the overall structure. The input activation is a *C*-dimensional FP16 vector. After the RMSNorm layer, we quantize the activation first from FP16 to INT4 using the standard Round-To-Nearest (RTN) method, then further transform it into four boolean variables. Next, the boolean activations are fed into a binarized fully-connected layers to compute K, Q, V matrices. The inner-loop computation is boolean, but the outputs recovers FP16 for query branch, and INT4 for key and value branches, *i.e.*, using 4-bit for KV cache. After the attention, we transform the activation vectors to boolean again so that all FC layer in the subsequent projection operations are also binarized.

At the core of quantization task lies the binarization of fully connected (FC) layers, which account for approximately 90% of memory bandwidth and computation. In contrast, other components such as normalization and attention score matrices contribute to the remaining 10%. Consequently, our work primarily focuses on the binarization of FC layers, which are utilized in computing key (K), query (Q), and value (V) matrices, as well as all projection layers. Since the most computationally expensive operation in these layers is the matrixvector multiplication within the inner loop, we pay most effort to design a computational efficiency FC layer, where the inner-loop only have binary operations.

Binarized fully-connected layer The core modification of the proposed BWA attention module is the binarization of the linear layer. To optimize the tradeoff between model accuracy and efficiency, we introduce a two-level weight grouping strategy for weight binarization, and a binarized decomposition to further quantize a 4-bit activations into four 1-bit boolean variables. Consequently, in the core inner-loop computation of the fully connected layers, only boolean operations are involved, drastically simplifying and boosting the computation.

Specifically, a FC layer computes a single token y = Wx, where $W \in \mathbb{R}^{C_{\text{in}} \times C_{\text{out}}}$, x is the C_{in} -dimensional vector represented as an input token, and y is the output token vector.

1) Channel-wise grouping Both weights and activations have a width dynamical range, which are hard to quantize using the same scaling parameter. We sort the input channels by the average activation scales of test samples X, *i.e.*, sort according to 242

243

244

245

246

247

196

197

198

199

200

201

202

203

204

205

206

207

297

298 299

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

326

295

diag(XX^T) in an ascending order. Then, following the sorted order, we group the channels with similar activation scale. In our implementation, we divide the $C_{out} = 4096$ channels into 32 groups. Each group has B = 128 channels.

248

249

254

257

262

264

266

267

269

271

272

273

276

277

279

281

285

293

$$y_{j} = \sum_{\ell=0}^{\lceil C_{\text{out}}/B \rceil - 1} \sum_{i=0}^{B-1} W_{j,(B\ell+i)} x_{B\ell+i}, \quad (1)$$
$$j = 1, 2, \dots, C_{\text{out}}$$

Note that channel-wise group is implemented by a proper permutation of rows and columns of the weight matrix. Therefore, channel-wise grouping does not introduce any dditional computational or memory cost.

2) Element-wise grouping Directly binarization of W to 1-bit is not enough resulting severe model degeneration forbidding practical usage, especially combined the target to reduce the activation to 4 bit or lower. Similar to BiLLM, we use additional element-level group to further group weight, and quantize them separately.

Define the set $D_{\ell,1} \subseteq \{0, 1, 2, \dots, B-1\}$, and its complementary $D_{\ell,0} = \{0, 1, 2, \dots, B-1\} \setminus D_{\ell,1}$. We quantized the two group separately using 1-bit for each weight element

$$W_{j,(B\ell+i)} \approx W_{j,(B\ell+i)} = \alpha_{j,\ell,s} q_{j,(B\ell+i)} + \beta_{j,\ell,s},$$
(2)

for s = 1 if $i \in D_{\ell}$ otherwise s=0, in which $\alpha_{j,\ell,s}$ and $\beta_{j,\ell,s}$ are two quantization parameters for each fine-grained group. In the standard LLAMA 7b model, each column-wise group contains has 128 channels. They are further categorized into two element-wise group, of which the sizes may not equal. The element-wise group are represented using a bit map. Together another bit $q_{j,i}$ representing weight sign. The weight matrix uses 1+1 bit per element.

This fine-grained level grouping can effectively the accuracy at the expense of increase additional computational cost. On the other hand, with proper usage of bitmap operations, the cost can be suppressed to be marginal. The element-wise grouping has been used in the works (Huang et al., 2024a; Dong et al., 2024), we improve this strategy by leveraging EM-based parameter searching with Hessian metric information to determine the optimal weighted split point.

3) Activation binarization. The input activation is quantized before feeding into the FC layer from FP16 to INT4 using the standard RTN method (3).

Specifically, the quantization process of RTN to obtain \mathbf{X}_q with *b* bits is expressed as

$$\mathbf{X}_q = \operatorname{clamp}(\lfloor \frac{\mathbf{X}}{\mu} \rceil + z, 0, 2^b - 1), \qquad (3)$$

where $\mu = \frac{\max(\mathbf{X}) - \min(\mathbf{X})}{2^{b} - 1}$ is the scaling parameter, $z = -\lfloor \frac{\min(\mathbf{X})}{\mu} \rceil$ is the shift parameter. The clamp function denotes restricting the quantization result to the range between $[0, 2^{b} - 1]$. The notion $\lfloor \rceil$ signifies the nearest rounding operation. μ is the quantization step size and z represents the zero point.

To further binarize the activation, we first transform the 4bit integer x_i to four 1bit boolean variables $b_{i,a}$ with a = 0, 1, 2, 3

$$x_i \approx \hat{x}_i = \mu p_i + z = \sum_{a=-1}^{3} \mu_a b_{i,a},$$
 (4)

where $\mu_a = 2^a \mu$ for a = 0, 1, 2, 3 and a special notation of $\mu_{-1} = z$ and $b_{i,-1} \equiv 1$ for the shift constant. Moreover, we can also relax μ_k as a free quantization parameter to be tuned manually or learn from data. Note that as the weight is reordered and grouped, the elements of the input activation vector will be permuted accordingly.

4) **Binarized FC layer** We utilize the commutative property of summation to make sure the inner multiplication summation is binary, so that the proposed method is actually boosting the computational speed.

Substituting (2) and (4)

$$y_j \approx \sum_{\ell=0}^{\lceil C_{\rm in}/B\rceil - 1} \sum_{i=0}^{B-1} \left[\left(\alpha_{j,\ell,s} q_{j,(B\ell+i)} + \beta_{j,\ell,s} \right) \right]$$
 32

$$\sum_{a=-1}^{3} \mu_{\ell,a} b_{(B\ell+i),a} \bigg]$$
322

$$=\sum_{\ell=0}^{\lceil \lceil C_{\rm in}/B \rceil - 1} \sum_{a=-1}^{3} \mu_{\ell,a} \sum_{s=0,1} 323$$

$$\left[\alpha_{j,\ell,s}v_{j,\ell,s} + \beta_{j,\ell,s}r_{j,\ell,s}\right],\tag{5}$$

with the bit-wise inner product $v_{j,\ell,s}$ and counting number of bits valued 1

$$v_{j,\ell,s} = \sum_{i \in D_{\ell,s}} q_{j,(B\ell+i)} b_{(B\ell+i)}, \ r_{j,\ell,s} = \sum_{i \in D_{\ell,s}} b_{(B\ell+i)}.$$
(6)

Since the above two summation only involves bit 328 variable $q_{j,(B\ell+i)}$, $b_{(B\ell+i)}$, and the set $D_{\ell,1}$ and 329 330 $D_{\ell,2}$ is represented by a bitmap, therefore, they can 331 be implemented efficiently by bit-wise XOR/AND 332 and pope operation.

333

335

336

337

338

341

345

351

354

355

367

370

373

For the bit-operation of (6), let $e = q \wedge b$. Then, v, r in (6) can be computed by

$$v_{j,\ell,s=0} = \operatorname{Popc} (e \wedge m),$$

$$v_{j,\ell,s=1} = \operatorname{Popc} (e \wedge (\neg m))$$

$$r_{j,\ell,s=0} = \operatorname{Popc} (b \wedge m),$$

$$r_{j,\ell,s=1} = \operatorname{Popc} (b \wedge (\neg m))$$
(7)

The inner loop, which originally be a multiplication-summation operation of 64 numbers, reduces to four bit-wise operations on 128 bits length variables, which can be efficiently implemented on both GPU and CPU.

5) Outlier activation In this paper, we trick the last channel-wise group as outlier, and use INT8 to quantize rather then the two bits following the work (Yuan et al., 2023; Zhao et al., 2024). It is hardware-friendly and can efficiency implemented via reordering the channels. Different from other methods, which usually use 256 channel for for outliers, we only use minimal 1 group, which only contain 128 channels for the outliers. This minimizes the outlier overhead to approximately 3% of the total channels in INT8 and representing all other normal channels in INT1. Experiments for set more outlier groups are done in Appendix, which can further improve the performance at the expense of costing extra bits.

3.2 Weight binarization and parameterization by Fine-Grained Group Hessian-Aware Quantization

In this subsection, we propose an EM-based method to determine the binarized values of the weight matrices and the associated parameters for dequantization.

The weight matrix is quantized according to (2), where $q_{j,(B\ell+i)}$ represents the binarized weight, and $\alpha_{j,\ell,s}$, $\beta_{j,\ell,s}$ denote the scaling and shifting parameters, respectively, used for dequantization. Furthermore, we also require a bit map matrix of identical dimension to the weight matrix, which serves to determine the fine-group affiliation of each corresponding element.

We start from minimizing the L2 norm of the weight matrix W, utilizing an approximate Hessian weight as proposed in (Hassibi and Stork,

1992; Frantar and Alistarh, 2022)

$$\mathcal{L}_{\hat{\mathbf{W}}} = \left\| \frac{1}{\operatorname{diag}(\mathbf{H}^{-1})} (\mathbf{W} - \hat{\mathbf{W}}) \right\|_{2}^{2}, \qquad (8)$$

where the Hessian matrix $H = XX^T$ is determined using a validation dataset, encapsulates information about the activation values during the computation of linear layers and gradient information during back propagation.

In our binarized parameterization (2), each channel (indexed by j) and each channel-wise group (indexed by ℓ) are parameterized independently. Without loss of generality, we only focus on a single channel-wise group, which has B weight elements $w_i, i \in 0, 1, ..., B - 1$. The binary representation of W with a binary element-wise group can at most have four different float-point values. Therefore, the quantization problem boils down to a 1-D clustering problem to determine 4 clusters, centered at $\hat{w}(0,0), \hat{w}(0,1), \hat{w}(1,0), \hat{w}(1,1)$. Formally, for each group, we solve the following minimization problem

$$\min_{\boldsymbol{s}, \boldsymbol{q} \in \{0,1\}^B, \hat{\boldsymbol{w}} \in \mathbb{R}^4} \sum_{i=0}^{B-1} \left(w_i - \hat{w}(s_i, q_i) \right)^2 / \text{diag}(\mathbf{H}^{-1})_i ,$$
(9)

where s, q represents the fine-group affiliation and the binary weight value respectively. Knowing the four values of \hat{w} , one can recover the scaling and shifting parameters $\alpha_{j,\ell,s}$ and $\beta_{j,\ell,s}$ by the method of undetermined coefficients formalized as a set of 4-D linear equations.

In practice, we spare the last channel-group for outlier activation using more bits, in which the corresponding weights are also quantized to the same type INT8. We implement a revised EM algorithm, as shown in Algorithm 1, to solve the above optimization problem with the last channel group as INT8 quantization.

3.3 Remarks on quantization strategies

We provide additional discussion on the quantization strategies employed in our approach, comparing them to existing methods such as the RTN quantization, vector quantization for weight quantization, and a recent work on binarized residual decomposition for activation quantization.

RTN quantization v.s. fine-group binary weight In our method, we utilize 1 bit to store binary weights and an additional 1 bit to represent element-wise group affiliation, effectively using 374

375

376

377

378

379

382

383

386

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

380 381

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

436

432

433

434

435

419

420

	K, outliers keep in IN18
	<i>iters</i> , EM steps
Ens	ure:
	B , weights after dequantization
	D , fine-grained group information
1:	$\mathbf{W} = \operatorname{reorder}(\mathbf{W}, \operatorname{diag}(\mathbf{X}\mathbf{X}^T))$
2:	$\mathbf{H} = 2\mathbf{X}\mathbf{X}^T$ > Hessian matrix
3:	$\mathbf{H}^{c} = \text{Cholesky}((\mathbf{H} + \lambda \mathbf{I})^{-1})$
4:	$\mathbf{B} = 0_{C_{in} \times C_{out}}$
5:	for $i = 0, B, 2B,, C_{out} - K - B$ do
6:	$\mathbf{W}^p = \mathbf{W}_{:,i:i+B}$
7:	$\mathbf{H}^{cp} = \operatorname{diag}(\mathbf{H}^{c}_{i:i+B \ i:i+B})$
8:	$\mathbf{C} = \text{init_centers}(\mathbf{W}^p, \mathbf{H}^{cp})$
9:	for $j = 1$ to <i>iters</i> do
10:	$\mathbf{D} = \text{get}_{groups}(\mathbf{W}^p, \mathbf{H}^{cp}, \mathbf{C}) \triangleright \text{E-step}$
11:	$\mathbf{I} = \text{get_clusters}(\mathbf{W}^p, \mathbf{H}^{cp}, \mathbf{C}, \mathbf{D})$
	M-step
12:	$\mathbf{C} = update_centers(\mathbf{W}^p, \mathbf{H}^{cp}, \mathbf{I}, \mathbf{D})$
	M-step
13:	end for
14:	$\mathbf{B}_{:,i:i+B} = \text{binary}(\mathbf{W}^p, \mathbf{C}, \mathbf{I}, \mathbf{D})$
15:	$\mathbf{E} = \frac{(\mathbf{W}^p - \mathbf{B})}{\mathbf{H}^{cp}}$
16:	$\mathbf{W}_{:i+B:Cout-K} = \mathbf{W}_{:i+B:Cout-K} - \mathbf{E}$
	$\mathbf{H}_{i\cdot i+B}^{c} \stackrel{i+B}{\underset{i+B}{\longrightarrow}} \mathbf{H}_{i\cdot i+B}^{c}$
17:	end for
18:	$\mathbf{B}_{:,Cout-K} = \text{quant_int8}(\mathbf{W}_{:,Cout-K})$
19:	Return B, D

Algorithm 1 Main Framework of our method

 $\mathbf{W} \in \mathbb{R}^{C_{out} \times C_{in}}$, weight matrix

-

 $\mathbf{X} \in \mathbb{R}^{T \times C_{in}}$, calibration data

B, block size

.1.

Require:

2 bits of information. It enables each weight element to take on four distinct values. In contrast to the RTN quantization method widely used in other post-training quantization (PTQ) methods, where dequantized values are equally spaced, our model allows the four values to be chosen arbitrarily, which is optimized by the proposed EM-based algorithm.

Vector quantization v.s. fine-group binary weight The optimization process for dequantization parameters in our approach is similar to that employed in vector-quantization-based methods (Frantar et al., 2022; van Baalen et al., 2024), where 2^n floating-point values are stored as representative values for *n*-bit quantization. In these approaches, dequantization must be performed before computing the vector inner product. In contrast, our method further parameterizes the representatives using binary weights and fine-grained group bits, along with floating-point scaling and shift parameters. This enables us to compute the vector inner product using pure Boolean operations as shown in (7), resulting in a significant boost in computational speed.

Binarized residual decomposition and 1×4 bit representation of activation The work (Zeng et al., 2024) explored the approach of transforming arbitrary integer weight and activation WxAa into $xa \times W1A1$ to achieve computational acceleration. It make use of bit operation to computation innerloop vector product, but the original work can not get good below W4A4. On the other hand, the expansion of high bits weight and activation usually result heavy over head, as the number of channel vectors (relates to memory bandwidth) roughly from (x + a) bits to (xa) bits. In this work, we manage to reduce $W(1 + 1)A(1 \times 4)$, together with bitmap operation on the fine-grain group, the over-head cost is marginal.

4 **Experiments**

Setup. We implemented our method on the Py-Torch (Paszke et al., 2019) framework, where all linear layer weights in the original model are quantized to 1+1 bit, and input activations of all linear layers are quantized to 1×4 bits. For weights, we adopt per-channel asymmetric quantization with a clipping ratio set to 1.0 across all experiments, utilizing the GPTQ quantization framework to compensate for quantization errors. For activations, we employ per-token asymmetric quantization with a clipping ratio of 1.0. To optimize performance, we use RTN for dynamic quantization of the activation matrix. For KV caches, we uniformly apply 4 bits quantization to store and load. The quantization group size is 128, and the number of outlier channels is 128 (approximately 3% of all channels). We use 128 random samples from the WikiText2 (Merity et al., 2016) training set as the calibration dataset, with a sequence length of 2048. The experiments related to the 7B model were conducted on a single NVIDIA GeForce RTX 3090 GPU, while the experiments involving the 13B model and acceleration efficiency were carried out on a single NVIDIA GeForce RTX A6000 GPU. All experiments were conducted more than three times, and the average values were recorded.

Models and Datasets. We apply our method

Table 1: Perplexity(\downarrow) and Zero-shot QA accuracy(\uparrow) results under the W4A4 and W2A4 settings on LLAMA1-7B and LLAMA2-7B. "FP16" denotes the performance of the original model represented in FLOAT16 format, with the best quantization performance highlighted in bold. The experimental results on the 13B model are presented in Table 3, Table 4, and Table 5.

Model	Bits	Method	Wiki.↓	РТВ↓	C4↓	PIQA↑	ARC-E↑	ARC-C↑	BoolQ↑	Hella.↑	Wino.↑	Avg.↑
	FP16	-	5.68	27.34	7.08	77.37	52.48	41.38	73.06	73.00	67.01	64.05
	WAAA	QuaRot	6.41	49.73	8.43	74.81	50.13	38.74	70.98	68.80	61.56	61.01
******	W+24+	Atom	6.30	30.28	7.98	75.35	51.60	36.69	70.86	67.27	64.33	62.21
LLAMAI -7B	W2AA	QuaRot	14.39	222.95	27.70	59.52	37.88	26.62	62.20	41.56	54.62	47.07
-78	W 2A4	Atom	16.65	298.78	33.87	57.24	35.23	26.11	53.98	36.77	50.51	43.31
	W(1+1)A16	BiLLM	35.04	421.27	39.59	61.20	36.00	25.70	62.70	36.80	51.10	45.58
	W(1+1)A4	BiLLM	18304	17152	20736	50.05	25.38	26.54	49.63	26.05	49.49	37.86
	W(1+1)A(1×4)	Ours	8.58	76.09	12.27	68.88	45.03	30.89	69.63	55.41	59.35	54.87
	FP16	-	5.47	22.51	6.97	76.93	53.58	40.53	71.07	72.96	67.17	63.71
	WAAA	QuaRot	6.32	71.21	8.67	74.32	51.60	38.23	68.41	69.24	61.56	60.89
11 4 4 4 4 2	W 4/44	Atom	6.18	27.94	8.05	75.24	52.74	37.12	71.16	67.89	63.93	62.58
-7B	W2AA	QuaRot	49.98	571.22	80.14	54.41	28.45	23.21	57.89	28.57	48.15	40.11
-78	112/14	Atom	19.49	508.82	39.85	56.69	32.32	23.21	58.53	35.74	49.49	42.66
	W(1+1)A16	BiLLM	32.48	3877.38	40.52	60.60	36.20	24.40	61.80	34.80	52.40	45.03
	W(1+1)A4	BiLLM	16128	17152	15168	50.22	26.30	27.90	45.23	26.10	49.88	37.61
	W(1+1)A(1×4)	Ours	8.89	69.46	12.74	68.72	46.13	30.55	66.12	55.76	58.01	54.22

to the open-source LLAMA1 (7B, 13B) (Touvron et al., 2023a), LLAMA2 (7B, 13B) (Touvron et al., 2023b), and Vicuna (7B, 13B) (Chiang et al., 2023) models and evaluate their performance on language generation and commonsense QA tasks. The primary metric for language generation tasks is perplexity, assessed on datasets including WikiText2, PTB (Marcus et al., 1994), and C4 (Raffel et al., 2020). For commonsense QA tasks, the main metric is zero-shot accuracy, evaluated on datasets such as PIQA (Bisk et al., 2020), ARC (Clark et al., 2018), BoolQ (Clark et al., 2019), HellaSwag (Zellers et al., 2019), and WinoGrande (Sakaguchi et al., 2021). Except for the C4 dataset, where we randomly select 256 samples of length 2048 from the test set for evaluation, we utilize the entire test set portion of these datasets for our testing.

Baseline. We compare our approach with stateof-the-art (SOTA) PTQ methods for weights and activations. Since few existing methods explore the W2A4 quantization setting, we implement W2A4 quantization for all compared methods to ensure fairness before evaluation. Our main baselines include Atom (Zhao et al., 2024), QuaRot (Ashkboos et al., 2024), and BiLLM (Huang et al., 2024a). Atom and QuaRot are SOTA methods under the W4A4 quantization setting, while BiLLM is the SOTA for the W(1+1)A16 quantization setting.

4.1 Main Results

486

487

488

489

490

491

492

493

494

496

497

498

499

502

505

507

510

511

512

513

514

Language Generation Tasks. We assess the perplexity of our method on language generation tasks
and conduct a fair comparison with existing SOTA
methods. As shown in Table 1 and Table 3, Atom

and QuaRot, as SOTA methods under the W4A4 setting, experience significant performance drops under the W2A4 setting. In contrast, our method significantly outperforms these methods on all datasets under the $W(1+1)A(1\times 4)$ setting which is equivalent to W2A4, and our method's perplexity evan approaches that of the FP16 model. It is noteworthy that BiLLM also utilizes an additional 1 bit to store extra fine-grained grouping information, thus we consider it as a W(1+1)A16 approach. As a similar method that employs fine-grained grouping like our method, its performance under the W(1+1)A16 configuration is significantly outperformed by our method with the $W(1+1)A(1\times 4)$ setting. Furthermore, when its activation values are quantized to 4 bits, the performance of BiLLM rapidly deteriorates.

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

Zero-Shot Tasks. We also evaluate our method on six important zero-shot tasks. Table 1 and Table 3 presents the comparison results between our method and the baselines. our method significantly outperforms existing methods under the W2A4 quantization setting and demonstrates stable accuracy, approaching the performance of the FP16 model.

4.2 Performance Analysis

Speedup. To evaluate the inference acceleration provided by our method, we adopt the kernel implementation from ABQ-LLM (Zeng et al., 2024), which supports decomposing arbitrary-dimensional WxAx operations into multiple W1A1 computations and leverages the acceleration effect of INT1 multiplication for significant speedup in matrix

610

611

612

613

614

615

616

617

618

619

584

585

586

587

588

589



Figure 3: A comparison of the time cost between the $W(1+1)A(1\times4)$ kernel and the INT4, INT8 kernels of CUTLASS for matrix multiplication on the A6000. More results are presented in Figure 4.

multiplication. We test the speedup of our method 552 553 under the $W(1+1)A(1\times4)$ setting on an A6000 compared to different bit-width settings supported by 554 CUTLASS, such as W8A8 and W4A4. Because the weights actually involved in the computations in our method are 1 bit, with an additional 1 bit 557 solely used for storing fine-grained grouping information, we consider our method as a quantization 559 method that can be viewed as a decomposition from 560 W2A4 downwards in terms of computational acceleration comparison. As shown in Figure 3, in terms of single-layer matrix computations, our method exhibits a more substantial speedup in comparison 564 to other bit-width settings, surpassing the kernel ac-565 celeration of CUTLASS by a factor of 3 in matrix 566 computations. This demonstrates that the approach of using INT1 for acceleration in our method can fully leverage the speedup benefits of low-bit com-569 putations. Moreover, since both the weights and activations in our method are quantized to very low 571 bit-widths, the additional computational overhead introduced by the decomposition does not significantly impact the gains achieved through INT1 574 computation.

4.3 Ablation Studies

576

577To evaluate the effectiveness of different quanti-578zation modules in our method, we compared the579accuracy gains or losses among various quantiza-580tion techniques employed within our method. The581results presented in Table 2 demonstrate that outlier582handling, Minimum distance quantization, and fine-583grained grouping, as the basic processing schemes

Table 2: Ablation experiments on the effects of different quantized components used in our method, with all experimental results based on LLAMA1-7B and a group size of 128.

Quantization Method	Wiki.↓
LLAMA-7B FP16	5.68
W1A4 GPTQ (Group size 128)	216713
+ Keep 128 outlier channels in INT8	6749
+ Minimum distance quantization	126.89
+ Fine-grained group, W(1+1)	8.69
+ Hessian-weighted distance metric	8.65
+ Binarized Residual Decomposition, A(1×4)	8.58

in our method. Each step significantly enhances the performance of the quantized model, effectively mitigating the performance collapse issue observed in the W1A4 quantized model. The introduction of the Hessian-weighted distance metric and binarized residual decomposition further boosts the quantized model's performance. Although numerically, the improvements in perplexity brought about by these two methods are not substantial, this is because the previous enhancements have already pushed the performance metrics close to those of the original model, leaving limited room for further improvement. Theoretically, the Hessian-weighted distance metric reveals a measure of weight importance, while the binarized residual decomposition elucidates the direction of performance enhancement after binarization decomposition.

5 Limitations

Although our our method can achieve fully binarized computation to reduce computational overhead during inference, a trade-off must be made regarding the accuracy of the quantized model. Consequently, we need to utilize additional bits to store the information of the quantized matrices, with the actual storage bits for weights and activations equivalent to 2 bits and 4 bits, respectively. This implies that our model has not been compressed to the theoretical extreme of a boolean model, leaving room for further improvement. Meanwhile, although the performance of our our method across various evaluation tasks is close to that of the prequantized model, it is not truly lossless quantization. This loss indicates that the quantized model has not fully restored the representational capacity of the original model. In the future, we consider employing methods that integrate quantization-aware

720

721

722

723

671

672

673

training to further enhance the efficiency and per-formance of our method.

References

627

629

631

642

643

667

- Saleh Ashkboos, Amirkeivan Mohtashami, Maximilian L Croci, Bo Li, Pashmina Cameron, Martin Jaggi, Dan Alistarh, Torsten Hoefler, and James Hensman. 2024. Quarot: Outlier-free 4-bit inference in rotated llms. arXiv preprint arXiv:2404.00456.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439.
 - Adrian Bulat, Yassine Ouali, and Georgios Tzimiropoulos. 2024. Qbb: Quantization with binary bases for llms. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
 - Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. *See https://vicuna. lmsys. org (accessed 14 April 2023)*, 2(3):6.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*.
- M Courbariaux and Y Bengio. 2016. Binarized neural networks: Training neural networks with weights and activations constrained to+ 1 or- 1. arxiv. *arXiv preprint arXiv:1602.02830*.
- Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. 2022. Gpt3. int8 (): 8-bit matrix multiplication for transformers at scale. *Advances in Neural Information Processing Systems*, 35:30318– 30332.
- Peijie Dong, Lujun Li, Yuedong Zhong, Dayou Du, Ruibo Fan, Yuhan Chen, Zhenheng Tang, Qiang Wang, Wei Xue, Yike Guo, et al. 2024. Stbllm: Breaking the 1-bit barrier with structured binary llms. arXiv preprint arXiv:2408.01803.
- Elias Frantar and Dan Alistarh. 2022. Optimal brain compression: A framework for accurate post-training quantization and pruning. *Advances in Neural Information Processing Systems*, 35:4475–4488.

- Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. 2022. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*.
- Babak Hassibi and David Stork. 1992. Second order derivatives for network pruning: Optimal brain surgeon. In *Advances in Neural Information Processing Systems*, volume 5. Morgan-Kaufmann.
- Wei Huang, Yangdong Liu, Haotong Qin, Ying Li, Shiming Zhang, Xianglong Liu, Michele Magno, and Xiaojuan Qi. 2024a. Billm: Pushing the limit of post-training quantization for llms. *arXiv preprint arXiv:2402.04291*.
- Wei Huang, Xudong Ma, Haotong Qin, Xingyu Zheng, Chengtao Lv, Hong Chen, Jie Luo, Xiaojuan Qi, Xianglong Liu, and Michele Magno. 2024b. How good are low-bit quantized llama3 models? an empirical study. *arXiv e-prints*, pages arXiv–2404.
- Haokun Lin, Haobo Xu, Yichen Wu, Jingzhi Cui, Yingtao Zhang, Linzhan Mou, Linqi Song, Zhenan Sun, and Ying Wei. 2024. Duquant: Distributing outliers via dual transformation makes stronger quantized llms. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Yifei Liu, Jicheng Wen, Yang Wang, Shengyu Ye, Li Lyna Zhang, Ting Cao, Cheng Li, and Mao Yang. 2024a. Vptq: Extreme low-bit vector post-training quantization for large language models. *arXiv preprint arXiv:2409.17066*.
- Zechun Liu, Changsheng Zhao, Igor Fedorov, Bilge Soran, Dhruv Choudhary, Raghuraman Krishnamoorthi, Vikas Chandra, Yuandong Tian, and Tijmen Blankevoort. 2024b. Spinquant–Ilm quantization with learned rotations. *arXiv preprint arXiv:2405.16406.*
- Liqun Ma, Mingjie Sun, and Zhiqiang Shen. 2024. Fbi-llm: Scaling up fully binarized llms from scratch via autoregressive distillation. *arXiv preprint arXiv:2407.07093*.
- Mitch Marcus, Grace Kim, Mary Ann Marcinkiewicz, Robert MacIntyre, Ann Bies, Mark Ferguson, Karen Katz, and Britta Schasberger. 1994. The penn treebank: Annotating predicate argument structure. In Human Language Technology: Proceedings of a Workshop held at Plainsboro, New Jersey, March 8-11, 1994.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*.
- Markus Nagel, Marios Fournarakis, Rana Ali Amjad, Yelysei Bondarenko, Mart Van Baalen, and Tijmen Blankevoort. 2021. A white paper on neural network quantization. *arXiv preprint arXiv:2106.08295*.

- 724 725 726 727
- 728 729 730

- 732 733 734
- 735 736 737 738 739
- 740 741 742
- 742 743 744
- 745 746 747 748 749 750
- 752 753 754 755 756 756

751

762

- 763 764 765 766
- 767 768
- 769 770 771
- 772 773 774

775 776 777

778

- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.
- Haotong Qin, Yifu Ding, Mingyuan Zhang, Qinghua Yan, Aishan Liu, Qingqing Dang, Ziwei Liu, and Xianglong Liu. 2022. Bibert: Accurate fully binarized bert. *arXiv preprint arXiv:2203.06390*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.
 - Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Mart van Baalen, Andrey Kuzmin, Markus Nagel, Peter Couperus, Cedric Bastoul, Eric Mahurin, Tijmen Blankevoort, and Paul Whatmough. 2024. Gptvq: The blessing of dimensionality for llm quantization. *arXiv preprint arXiv:2402.15319*.
- Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Huaijie Wang, Lingxiao Ma, Fan Yang, Ruiping Wang, Yi Wu, and Furu Wei. 2023. Bitnet: Scaling 1-bit transformers for large language models. *arXiv preprint arXiv:2310.11453*.
- Hongyu Wang, Shuming Ma, and Furu Wei. 2024. Bitnet a4. 8: 4-bit activations for 1-bit llms. arXiv preprint arXiv:2411.04965.
- Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. 2023. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*, pages 38087–38099. PMLR.
- Yuzhuang Xu, Xu Han, Zonghan Yang, Shuo Wang, Qingfu Zhu, Zhiyuan Liu, Weidong Liu, and Wanxiang Che. 2024. Onebit: Towards extremely low-bit large language models. *arXiv preprint arXiv:2402.11295*.

Zhihang Yuan, Lin Niu, Jiawei Liu, Wenyu Liu, Xinggang Wang, Yuzhang Shang, Guangyu Sun, Qiang Wu, Jiaxiang Wu, and Bingzhe Wu. 2023. Rptq: Reorder-based post-training quantization for large language models. *arXiv preprint arXiv:2304.01089*. 779

780

781

783

785

787

788

789

790

791

792

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*.
- Chao Zeng, Songwei Liu, Yusheng Xie, Hong Liu, Xiaojian Wang, Miao Wei, Shu Yang, Fangmin Chen, and Xing Mei. 2024. Abq-llm: Arbitrary-bit quantized inference acceleration for large language models. *arXiv preprint arXiv:2408.08554*.
- Yilong Zhao, Chien-Yu Lin, Kan Zhu, Zihao Ye, Lequn Chen, Size Zheng, Luis Ceze, Arvind Krishnamurthy, Tianqi Chen, and Baris Kasikci. 2024. Atom: Lowbit quantization for efficient and accurate llm serving. *Proceedings of Machine Learning and Systems*, 6:196–209.

A Additional Experimental Results

Time required for quantization. Our method quantizes the weight matrices within all linear layers of the full-precision model. The quantization process for the 7B model can be completed in approximately 20 minutes, while the 13B model requires only about 30 minutes.

Results on 13B models. As shown in Table 4 and Table 5, we evaluated the performance of our method and other quantization methods on language generation tasks and zero-shot QA task accuracy using LLAMA1-13B and LLAMA2-13B models. Our findings indicate that, in general, the model performance adheres to the principle that increasing the number of model parameters leads to improved model performance. Furthermore, our method achieved state-of-the-art results across all evaluated metrics.

Results of different outlier channel number settings. In Table 6, we compare the relationship between different numbers of outlier channels and the quantization performance of our method. Since the group size is set to 128, we also use 128 as the unit here. The results demonstrate that preserving a small number of outliers with high precision can ensure overall quantization performance. Furthermore, when the number of outlier channels is increased, the model performance exhibits a nearly linear upward trend, with only a modest overall improvement. Therefore, we adopt 128 outlier channels as our baseline setting.

Table 3: Perplexity(\downarrow) and Zero-shot QA accuracy(\uparrow) results under the W4A4 and W2A4 settings on Vicuna family. "FP16" denotes the performance of the original model represented in FLOAT16 format, with the best quantization performance highlighted in bold.

Model	Bits	Method	Wiki.↓	PTB↓	C4↓	PIQA↑	ARC-E↑	ARC-C↑	BoolQ↑	Hella.↑	Wino.↑	Avg.↑
	FP16	-	6.78	26.78	8.55	77.80	56.06	39.93	75.69	71.06	67.80	64.72
	WAAA	QuaRot	7.80	52.44	10.87	73.67	53.20	37.71	72.45	67.66	60.93	62.12
Vicuna	w4A4	Atom	7.22	31.75	9.36	75.14	55.60	37.63	77.25	67.08	64.40	64.42
-v1.5-7B	-v1.5-7B	QuaRot	39.51	226.50	65.17	55.66	33.38	22.75	62.08	31.71	50.51	44.03
	WZA4	Atom	15.96	107.68	25.13	56.64	31.90	29.61	64.07	46.30	55.33	47.31
	W(1+1)A(1x4)	Ours	9.51	45.61	13.35	69.97	50.00	33.45	71.96	57.81	59.43	57.10
	FP16	-	5.95	25.15	7.78	78.40	56.44	44.80	76.51	74.63	69.06	66.64
	WAAA	QuaRot	6.81	54.16	9.64	74.81	51.43	40.53	70.73	70.96	62.12	62.51
Vicuna	w4A4	Atom	6.32	27.64	8.25	76.44	54.67	43.34	74.83	72.07	66.46	65.37
-v1.5-13B	W2AA	QuaRot	18.32	273.86	37.86	56.69	36.49	26.11	62.42	38.40	53.59	45.54
	W 2/44	Atom	19.84	174.63	36.39	54.95	34.13	25.68	61.74	37.55	52.17	44.37
	W(1+1)A(1×4)	Ours	7.91	49.71	11.45	71.44	52.36	38.65	68.93	62.34	62.35	59.35

Table 4: Perplexity(\downarrow) results under the W4A4 and W2A4 settings on LLAMA1-13B and LLAMA2-13B. "FP16" denotes the performance of the original model represented in FLOAT16 format, with the best quantization performance highlighted in bold.

Model	Dite	Method	Perplexity↓			Model	Bite	Method	Perplexity↓		
Widdei	Dits		Wiki.	РТВ	C4	Niouei	Dits	Methou	Wiki.	РТВ	C4
	FP16	-	5.09	19.23	6.61		FP16	-	4.88	28.87	6.47
	WAAA	QuaRot	5.71	36.10	7.57		W4A4	QuaRot	5.59	64.27	7.84
LLAMA1	W4/14	Atom	5.47	22.16	7.04	LLAMA2		Atom	5.26	32.46	6.95
-13B	W2 A 4	QuaRot	11.14	156.30	20.80	-13B	W2 A 4	QuaRot	17.49	386.40	38.88
	W2A+	Atom	11.69	115.62	19.55		W 2/14	Atom	11.24	152.68	18.15
	W(1+1)A(1×4)	Ours	7.19	37.20	10.18		W(1+1)A(1×4)	Ours	7.17	56.91	10.44

Table 5: Zero-shot QA accuracy(\uparrow) results under the W4A4 and W2A4 settings on LLAMA1-13B and LLAMA2-13B. "FP16" denotes the performance of the original model represented in FLOAT16 format, with the best quantization performance highlighted in bold.

Model	Dite	Mathad	athod Zero-shot Accuracy↑									
Widdei	Bits	Methou	PIQA	ARC-E	ARC-C	BoolQ	HellaSwag	WinoGrande	Avg.			
	FP16	-	79.05	59.89	44.71	68.47	76.23	70.24	66.43			
	W/AAA	QuaRot	76.61	55.30	41.64	67.09	73.02	65.59	64.24			
T T AMA 1 12D	W4A4	Atom	77.64	58.38	41.81	68.50	73.75	65.98	64.64			
LLAMAI-13B	W2A4	QuaRot	64.42	41.50	28.75	63.36	48.49	56.67	50.53			
	W 2A4	Atom	59.68	36.62	29.01	58.56	44.84	52.01	46.79			
	W(1+1)A(1×4)	Ours	72.09	48.57	34.13	62.54	62.63	64.88	57.47			
-	FP16	-	79.00	57.95	44.28	69.02	76.58	69.69	66.09			
		QuaRot	76.93	52.10	40.70	68.10	72.70	62.51	62.99			
LLAMA2-13B	W 4/14	Atom	77.37	56.73	42.32	67.62	74.07	68.27	65.46			
	W2A4	QuaRot	59.09	34.60	24.23	62.11	35.03	51.38	44.41			
		Atom	61.15	40.36	29.52	61.56	45.62	51.22	48.24			
	W(1+1)A(1×4)	Ours	71.98	49.92	36.26	65.90	60.52	61.80	57.73			

Comparison of different kernels. In Figure 4, we comprehensively evaluate the performance of the $W(1+1)A(1\times4)$ kernel and the INT8, INT4 kernels from CUTLASS, based on the matrix multiplication sizes that may occur in the LLAMA model. Since our method incorporates a small amount of INT8 mixed-precision quantization, for the handling of outliers, we separately measure the computational efficiency of outliers and normal values. Subsequently, we derive the overall computational efficiency by considering the proportion of these

two components.

Model Size. We present in Table 7 the theoretical compression effectiveness of our method on models of various sizes within the LLAMA family. In our calculation of the model size, we have included both the quantization parameters and the additional storage incurred by fine-grained grouping, which results in our findings being slightly larger than those reported in BiLLM (Huang et al., 2024a). The binarization of weights significantly reduces the storage size of quantized LLMs and the GPU

Table 6: The impact of different outlier channel number settings of the quantized model on the perplexity (\downarrow) and the zero-shot QA accuracy(\uparrow). "FP16" denotes the performance of the original model represented in FLOAT16 format.

Model	Ch.	Wiki.↓	PTB↓	C4↓	PIQA↑	ARC-E↑	ARC-C↑	BoolQ↑	Hella.↑	Wino.↑	Avg.↑
	FP16	5.68	27.34	7.08	77.37	52.48	41.38	73.06	72.99	67.01	64.05
	0	471.19	1025.28	228.17	53.59	28.75	24.57	50.73	28.13	50.51	39.38
*******	128	8.58	76.09	12.27	68.88	45.03	30.89	69.63	55.41	59.35	54.87
LLAMAI 7D	256	8.20	65.97	11.70	69.75	45.33	32.42	65.87	56.31	57.85	54.59
-70	512	7.80	57.44	10.90	71.27	47.94	34.22	65.57	58.46	59.19	56.11
	768	7.52	52.06	10.44	71.38	47.31	34.04	66.12	60.09	61.17	56.69
	1024	7.26	50.42	9.95	72.14	47.01	34.39	69.30	61.16	61.01	57.50
	FP16	5.47	22.51	6.97	76.93	53.58	40.53	71.07	72.96	67.17	63.71
LLAMA2	128	8.89	69.46	12.74	68.72	46.13	30.55	66.12	55.76	58.01	54.22
-7B	256	8.52	61.01	12.16	69.97	47.64	31.57	68.50	56.22	59.19	55.51
	512	8.00	56.77	11.43	69.80	46.97	31.66	67.83	57.55	60.77	55.76
	FP16	6.78	26.78	8.55	77.80	56.06	39.93	75.69	71.06	67.80	64.72
Vicuna	128	9.51	45.61	13.35	69.97	50.00	33.45	71.96	57.81	59.43	57.10
-v1.5-7B	256	9.28	44.01	12.94	70.35	51.05	34.04	72.97	58.42	61.40	58.04
	512	8.88	41.49	12.42	71.87	50.76	34.22	73.79	58.80	64.09	58.92



Figure 4: A comparison of the computational efficiency between the $W(1+1)A(1\times4)$ kernel and the INT4, INT8 kernels of CUTLASS for matrix multiplication with varying input lengths is conducted on the A6000.

memory and bandwidth requirements during inference. Across LLAMA models of different sizes, our method achieves a compression ratio of over 5×.

Table 7: Model size comparison of LLAMA family.

Models	FP16	Ours
LLAMA-7B	13.5GB	2.69GB
LLAMA-13B	24.2GB	4.82GB
LLAMA-30B	60.5GB	12.05GB
LLAMA-65B	121.0GB	24.11GB

B Ai Assistants in Research or Writing

We use a local LLama3.3 model to polish the draft for checking grammar and improving expression. Research ideas, experiment design and discussion contents are all original by the authors.