ClusComp: A Simple Paradigm for Model Compression and Efficient Finetuning

Anonymous ACL submission

Abstract

As large language models (LLMs) scale, model compression is crucial for edge deployment and accessibility. Weight-only quantization reduces model size but suffers from performance degradation at lower bit widths. Moreover, standard finetuning is incompatible with quantized models, and alternative methods often fall short of full finetuning. In this paper, we propose ClusComp, a simple yet effective compression paradigm that clusters weight matrices into codebooks and finetunes them block-by-block. ClusComp (1) achieves superior performance in 2-4 bit quantization, (2) pushes compression to 1-bit while outperforming ultra-low-bit methods with minimal finetuning, and (3) enables efficient finetuning, even surpassing existing quantization-based approaches and rivaling full FP16 finetuning. Notably, ClusComp supports compression and finetuning of 70B LLMs on a single A6000-48GB GPU.

1 Introduction

001

002

004

011

013

017

021

022

037

041

Large language models (LLMs) have garnered significant acclaim and success across various domains and applications (Dubey et al., 2024; Brown et al., 2020; Raffel et al., 2020a). With ongoing advancements, the scope and complexity of released LLMs have witnessed exponential growth, with some LLMs encompassing >50B parameters (DeepSeek-AI et al., 2024; Zhang et al., 2022; Scao et al., 2022). This remarkable upscaling introduces considerable challenges, particularly when deploying these models or granting their accessibility to users with constrained resources. To address these challenges, weight-only post-training quantization (PTQ) has emerged as an effective approach.

PTQ methods can generally be classified into three categories: statistic-based, gradient-based, and codebook-based approaches. Statistic-based methods (Dettmers et al., 2024; Lin et al., 2024; Frantar et al., 2022) determine the quantization grid based on the weight value distribution, whereas gradient-based methods (Shao et al., 2024; Ma et al., 2024) optimize the quantization grid with some calibration samples. Codebook-based methods (Egiazarian et al., 2024; van Baalen et al., 2024; Kim et al., 2024; Park et al., 2024) cluster similar weight elements to the shared quantized centroids, employing non-uniform quantization and pushing the limits to extremely low bit levels. However, these methods continue to struggle with low-bit quantization and the presence of outliers in weights, leading to significant performance degradation. 042

043

044

047

048

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

078

079

081

Another challenge PTQs encounter is their limited support for finetuning, which is crucial for adapting LLMs to various downstream tasks. Finetuning LLMs is computationally expensive due to their large scale and the need to cache activations and store optimizer states. PTQ, which compresses LLMs, appears to be a promising approach for finetuning as it reduces the memory requirements for loading LLMs. However, most quantization techniques use a round-to-nearest operation, which does not support gradient backpropagation. Typically, parameter-efficient methods (Dettmers et al., 2023; Li et al., 2024c; Liao and Monz, 2024b) are employed to train the added parameters while keeping the quantized LLMs frozen, bypassing this limitation. However, this fine-tuning approach presents two major drawbacks: (1) freezing quantized LLMs prevents further reduction of quantization errors during finetuning; (2) The low-rank nature of most parameter-efficient methods restricts their expressiveness (Biderman et al., 2024; Liao and Monz, 2024a).

In this paper, we propose a simple while effective paradigm that mainly applies **Clus**tering to **Comp**ress LLMs, referred to as *ClusComp*. Additionally, ClusComp can function as a parameter and memory-efficient finetuning method. Our preliminary experiments reveal that LLMs are increasingly difficult to quantize, mainly due to the growing fre-



Figure 1: Left: Compression quality of ClusComp. Only the most competitive baselines are shown here. Please refer to Table A.1 and A.2 for a full comparison. Methods in triangle use more calibration samples than the ones in circle. **Right:** Compression efficiency of ClusComp that consists of sequential clustering and reconstruction steps. E.g., compressing a 70B LLM takes 25.6h + 2GB for clustering, and 3.5h + 26GB for reconstruction on 1 GPU.

110

111

112

113

114

quency of outliers in their weight matrices (§3). Based on this observation, we propose using clustering instead of quantization to compress LLMs, retaining all values in FP16 format to circumvent issues arising from outlier quantization ($\S3.1$). To further reduce compression errors, we minimize block-wise output discrepancies between the compressed and uncompressed blocks, using a limited set of calibration samples $(\S3.3)$. In addition, by incorporating an inexpensive end-to-end recovery finetuning step, we can even push the compression rate to the 1-bit level. Since all parameters remain in FP16 after compression, ClusComp fully supports standard neural training, thus allowing the seamless finetuning of compressed LLMs on various downstream tasks (§3.4).

We begin by evaluating the effectiveness of Clus-Comp in the context of model compression across 2 language modeling tasks and 6 zero-shot reasoning tasks. ClusComp consistently surpasses various baselines at 2-4 levels, even achieving a perplexity of <13 at the 2-bit level on WikiText2 (Merity et al., 2017) for all LLMs (§4.1). Following recovery finetuning, ClusComp's performance at 2-bit and 1-bit levels approaches that of the FP16 model, with an accuracy of 57.8 vs 68.6 for the 2-bit Llama-3-8B and 51.4 vs 75.4 for the 1-bit Llama-3-70B (§4.2). Additionally, ClusComp demonstrates its utility as a parameterefficient (< 1%) and memory-efficient (42GB for Llama-3-70B) finetuning method, outperforming quantization-based and memory-efficient finetuning approaches, while matching the performance of full finetuning (§4.3).

115

116

117

118

119

120

121

122

123

125

126

127

128

129

130

131

132

134

135

136

137

138

139

140

141

142

2 Related Works

Quantization refers to the process of converting floating-point values into discrete levels, thereby reducing the bit-width required and minimizing memory consumption during model loading. Taking the symmetric uniform quantization as an example, a weight matrix **W** is quantized as follows:

$$W_q = \operatorname{clamp}(\lfloor \frac{W}{s} \rceil, -2^{b-1}, 2^{b-1} - 1)$$
 124

where the scale factor $s = \frac{\max(|W_{\min}|, |W_{\max}|)}{2^{b}-1}$, *b* denotes the bit-width, and [] represents the round-tonearest (RTN) operation. Since the quantization grid is uniform, its effectiveness is contingent on the distribution of the weight values. In cases where the weight matrix contains a significant number of outliers or is quantized to lower bit-widths, the resulting quantization error may be substantial.

Post-training quantization (PTQ) methods, such as GPTQ (Frantar et al., 2022), AWQ (Lin et al., 2024), and OmniQuant (Shao et al., 2024), apply quantization to a model after training with minimal computational resources. However, these approaches, which rely on uniform quantization, are significantly impacted by the presence of outliers in the weight matrices. Recent methods (Dettmers et al., 2024; Yuan et al., 2024; Huang et al., 2024a) address this challenge by retain-



Figure 2: The Llama series becomes progressively harder to quantize. Left & Middle: From Llama-2 to Llama-3, all methods struggle more with lower-bit quantization, though ClusComp⁻ is the least affected. Right: From Llama-1 to Llama-3, the average weight kurtosis across most layers increases.

ing salient weights in FP16 format, thereby main-143 taining strong performance at lower bit widths. 144 Nonetheless, these mixed-precision quantization 145 techniques require specially optimized CUDA ker-146 nels to either enhance or preserve inference speed. 147 Closely related to our proposed method, ClusComp, 148 are works such as GPTVQ (van Baalen et al., 2024), 149 QuIP# (Tseng et al., 2024) and SqueezeLLM (Kim et al., 2024) which implement quantized codebooks 151 for non-uniform quantization, achieving state-of-152 the-art performance for ultra-low-bit quantization. 153 ClusComp, however, differs in two significant 154 ways: (1) The codebook in ClusComp is stored 155 in FP16, offering additional advantages for subse-156 quent recovery training and finetuning; (2) While 157 other methods face limitations similar to those in 158 VAE-like approaches (Kingma and Welling, 2014), 159 where large and high-dimensional codebooks are 160 infeasible due to mode collapse, ClusComp circumvents this issue. Our fixed-code design allows us to utilize a codebook size of 2^{16} in 4-16D without 163 encountering such difficulty. 164

Finetuning is crucial for adapting LLMs to various domains and applications. Quantization, which reduces model size, is theoretically more conducive to finetuning. However, directly finetuning a quantized model is not a standard approach, as RTN does not support gradient backpropagation. Finetuning using a straight-through estimator (STE) (Bengio et al., 2013) is relatively underexplored and may lead to catastrophic forgetting (Malinovskii et al., 2024). Previous works (Xu et al., 2024a; Liao and Monz, 2024b; Dettmers et al., 2023) propose freezing the quantized model while updating newly added LoRAs (Hu et al., 2022). However, these approaches suffer from two key limitations: (1) Freezing the quantized model prevents the mitigation of quantization errors during finetuning. Moreover, not all quantization meth-

165

166

168

170

171

172

173

174

175

177

178

179

181

ods are suitable for finetuning. Popular techniques like GPTQ and QLoRA, while widely used, exhibit significant quantization errors below 4-bit. (2) The expressiveness of LoRA is constrained by its bottleneck design (Biderman et al., 2024; Liao and Monz, 2024a). In contrast, ClusComp, where all parameters are maintained in high precision, inherently supports seamless finetuning. Additionally, updating the codebook in ClusComp results in modifying all parameters in the weight matrices, even offering superior performance compared to full finetuning while maintaining a similar number of trainable parameters as LoRA. 182

183

184

185

186

187

188

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

208

209

210

211

212

213

214

215

216

217

218

3 ClusComp

Pilot study. Before introducing ClusComp, we present a key observation from our experiments on different Llama series. As in Figure 2 (Left & Middle), when reducing the bit-width, the performance of various quantization methods (RTN, GPTQ, and AWQ) follows a similar trend: Llama-3 (Dubey et al., 2024) proves more challenging to quantize than Llama-2 (Touvron et al., 2023b).

We hypothesize that the increasing quantization difficulty stems from a higher frequency of outliers in Llama-3's linear layers. Since these methods rely on uniform quantization, they are particularly sensitive to outliers in weight matrices. To test this, we analyzed the kurtosis of weight matrices—a well-established metric for detecting outliers (Bondarenko et al., 2023). As shown in Figure 2 (Right), we observe: (1) Higher kurtosis at the beginning and end of all models; (2) A consistent increase in kurtosis from Llama-1 to Llama-3 across most layers, indicating more frequent outliers. This trend likely explains the growing quantization challenge and suggests that future Llama models may be even harder to quantize.¹

¹Further discussion and proof are in §A.1.

Since quantization performance is affected by outliers, could an alternative compression approach involve storing all weight values in FP16 rather than applying quantization?

The first idea that comes to mind is clustering, where similar weight values are represented by a single shared value. This approach enables model compression while maintaining all weight values in FP16. In the following, we introduce three ClusComp variants that leverage clustering for LLM compression: *ClusComp*⁻, which applies clustering alone; *ClusComp*, which enhances ClusComp⁻ with block-wise error minimization; and *ClusComp*⁺, which further refines the compressed model through next-token prediction.

3.1 ClusComp⁻: Clustering

219

220

228

232

233

234

236

240

241

242

243

244

245

246

247

248

251

253

258

260

Considering a weight matrix $W \in \mathbb{R}^{d_{in} \times d_{out}}$, direct clustering along either dimension of W is suboptimal as it leads to a significant reconstruction error, particularly due to the large values of d_{in} and d_{out} in LLMs. To mitigate this issue, we reshape Winto a set of lower-dimensional vectors, denoted as $W' = \{w_1, w_2, \dots, w_k\}$, where each $w_i \in \mathbb{R}^g$ and $k = \frac{d_{in} \cdot d_{out}}{g}$.² The goal is to partition W' into n clusters $\{C_1, C_2, \dots, C_n\}$ by solving the following optimization problem:

$$\operatorname{argmin}_{\{C_1, C_2, ..., C_n\}} \sum_{j=1}^n \sum_{\boldsymbol{w}_i \in C_j} ||\boldsymbol{w}_i - \boldsymbol{c}_j||^2$$

where $c_j \in \mathbb{R}^g$ denotes the centroid of cluster C_j . This clustering problem is well-established in the machine learning literature and can be iteratively addressed using K-means (Lloyd, 1982) with the Expectation-Maximization (EM) algorithm:

- E-step: Each vector w_i is assigned to the cluster whose centroid c_j minimizes the Euclidean distance, i.e., $C_j^{(t)} = \{w_i : ||w_i c_j^{(t)}||^2 \leq ||w_i c_l^{(t)}||^2 \quad \forall l\}.$
- **M-step**: The centroid of each cluster is updated as the mean of the vectors assigned to that cluster, i.e., $c_j^{(t+1)} = (\sum_{w_i \in C_j^{(t)}} w_i) / |C_j^{(t)}|.$

Upon completion, two key elements are obtained for each weight matrix: (1) a codebook

Setting	#Params. for Codes	#Params. for Codebook	\bar{b}
g4n65500	4.19M	0.26M (1.55%)	4.25
g6n65500	2.80M	0.39M (2.32%)	3.04
g9n65500	1.86M	0.59M (3.52%)	2.34

Table 1: Bits for **W** with $d_{in}, d_{out} = 4096$ (16.78M).

 $C = \{c_1, c_2, \ldots, c_n\} \in \mathbb{R}^{g \times n}$ that contains all centroids, and (2) a set of codes $q = \{q_1, q_2, \ldots, q_k\} \in \{1, 2, \ldots, n\}^k$ that records the assignment of each vector w_i to the closest centroid, where $q_i = q(w_i) = j$ if $c_j = argmin_{c_l \in C} ||w_i - c_l||^2$. Using the codes q and the codebook C, the weight matrix W' can be reconstructed as $\hat{W}' = \{c_{q_1}, c_{q_2}, \ldots, c_{q_k}\}$. In PyTorch (Paszke et al., 2017), the linear layer is adapted as in Listing 1.

Remark: ClusComp, when applied solely with clustering, is referred to as ClusComp⁻. In this configuration, only the weight matrices (no calibration samples) are utilized, leading to substantial memory efficiency, with a mere 2GB memory consumption for both 7B and 70B LLM on 1 GPU, as shown in Figure 1. Moreover, this process can be considerably accelerated with more GPUs, as the clustering of different matrices is independent.

3.2 Estimate model size

After clustering, it is sufficient to store the codes $q \in \{1, 2, ..., n\}^k$ and codebook $C \in \mathbb{R}^{g \times n}$. Unlike prior works (van Baalen et al., 2024; Egiazarian et al., 2024; Tseng et al., 2024), we don't quantize the codebook; instead, we store it in the FP16 format. The bit-width required for the codes depends on the range of n. To maintain efficiency, we set $n < 2^{16}$ and use unsigned 16-bit integers to represent the codes. Thus, the average bits-perparameter is calculated as:

$$\bar{b} = \frac{\text{size in bits}}{\text{number of parameters}}$$

$$= \frac{16 \cdot k + 16 \cdot g \cdot n}{d_{\text{in}} \cdot d_{\text{out}}} = \frac{16}{g} + \frac{16 \cdot g \cdot n}{d_{\text{in}} \cdot d_{\text{out}}}$$
(1)
292

In the right-most of Equation (1), the first term corresponds to the bit-width allocated to the codes, and the second term corresponds to the bit-width of the codebook. E.g., for a linear layer W with $d_{\rm in} = d_{\rm out} = 4096$, clustering with g = 4 and $n = 2^{16} - 1$ results in $\bar{b} \approx 4 + 0.25 = 4.25$. This demonstrates that the majority of the bit-width is allocated to the codes, which is a primary reason for constraining $n < 2^{16}$, so we can use 16 instead 290

293

294

295

296

297

299

²In cases where the dimensions are not divisible, zeropadding is applied to W. In PyTorch, the matrix is reshaped as W' = W.transpose(1, 0).view(-1, g), where transposing W offers slightly better performance.

307

311

312

313

314

315

316

317

319

322

323

332

333

334

337

338

339

342

347

350

302

of 32-bit integers to represent the code. Further reducing n to a smaller range leads to fewer centroids, which in turn increases reconstruction error. More settings can be found in Table 1 and A.1.

Remark: While we express the model size in terms of bits-per-parameter, it is important to note that no quantization is applied in ClusComp. Instead, we reduce the number of parameters in W from $d_{in} \cdot d_{out}$ to $k + g \cdot n$. Utilizing bits-per-parameter allows for a direct comparison between ClusComp and quantization-based methods. Surprisingly, ClusComp⁻, which only incorporates the clustering step without employing any calibration data, already surpasses RTN, GPTQ and AWQ, as illustrated in Figure 2 (Left & Middle).

3.3 ClusComp: Block error minimization

Block-wise error minimization (block-wise reconstruction or knowledge distillation) has emerged as a standard, efficient and effective approach to reducing quantization error (Egiazarian et al., 2024; Tseng et al., 2024; van Baalen et al., 2024; Shao et al., 2024; Liao and Monz, 2024b). To further mitigate the compression error caused by clustering, we incorporate block-wise error minimization into ClusComp⁻ using a limited set of calibration samples, expressed as:

 $\operatorname{argmin}_{C_s} ||\mathcal{F}(Ws, X) - \mathcal{F}(Cs, qs, X')|| \quad (2)$

Here, \mathcal{F} denotes a Transformer block (Vaswani et al., 2017), Ws represent the weight matrices in the uncompressed block, and Cs and qs denote the codebooks and codes in the compressed block. X refers to the input of the uncompressed block, which is also the output from the previous uncompressed block, while X' is the input to the compressed block, originating from the output of the preceding compressed block. For the first block, we have X = X'. Block-wise error minimization is memory-efficient as it only requires loading two blocks into the GPU simultaneously.

Remark: In Equation (2), we only train the codebook Cs while keeping the codes qs fixed as indices. This design offers two key advantages: (1) It enhances data efficiency. As illustrated in Table 1, the majority of parameters are represented by the codes. Training both the codebooks and codes with a limited number of calibration samples (128) leads to overfitting; (2) More importantly, training the codes with a large number of centroids (2¹⁶) can result in mode collapse (Sun et al., 2024b; Kingma and Welling, 2014). Since the codes already exhibit a uniform distribution after clustering (see Figure A.2), keeping the codes fixed indicates that all centroids in the codebook can be trained uniformly. Such a code-fixed design is also applied to the following recovery and finetuning step. Combining both clustering and block-wise error minimization steps, we term this method ClusComp. 351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

379

380

381

382

383

384

385

386

388

389

390

391

392

393

394

395

396

397

398

399

3.4 ClusComp⁺: Recovery and finetuning

We present the adapted linear layer for ClusComp in Listing 1, which can be seamlessly integrated as a replacement for the original linear layer in LLMs. As the codebook is represented in FP16, this new layer inherently supports training without requiring additional tricks, like STE.

Recovery training. The compressed LLMs can be further trained by predicting the next token to recover information lost due to compression. This is achieved by finetuning the codebook parameters. This form of training is memory-efficient in two distinct ways: (1) Since the LLM is already compressed, loading it onto the GPU consumes less memory compared to the FP16 version; (2) As illustrated in Table 1, the parameters in the codebook account for < 5% of the total parameters in the FP16 version, making the training both parameterefficient and memory-efficient. We refer to Clus-Comp with recovery training as ClusComp⁺.

Finetuning. Like recovery training, finetuning the compressed LLM on downstream tasks can also be performed efficiently. Unlike QLoRA (Dettmers et al., 2023), which freezes the quantized LLM and trains only the LoRA (Hu et al., 2022) modules, finetuning the codebook alone eliminates the need for this additional constraint. This approach offers two key advantages over QLoRA: (1) Freezing the quantized LLM prevents mitigation of quantization errors, whereas finetuning the codebook can further address compression errors for downstream tasks; (2) The low-rank bottleneck of LoRA limits its expressiveness (Biderman et al., 2024). In contrast, finetuning the codebook is analogous to adapting the entire high-rank weight matrix, providing greater flexibility and expressiveness.

4 **Experiments**

4.1 Compression results

LLMs and evaluation. We evaluate ClusComp on widely adopted LLM families: Llama-1-7B, Llama-2-7B/13B/70B and Llama-3-8B/70B (Tou-



Figure 3: Average zero-shot accuracy over 5 or 6 commonsense reasoning tasks, only including competitive baselines. Please refer to Table A.3 and A.4 for detailed numbers and the full comparison.

Method	#Bit	AI2D ↑	ChartQA ↑	DocVQA ↑	$\textbf{MMBench} \uparrow$	Avg ↑	$\mathbf{MME}~(\mathrm{cog}~/~\mathrm{per})\uparrow$
LLaVA-Next-8B	16.00	71.7	69.2	78.2	72.2	72.8	1965.1 (376.8 / 1588.3)
GPTQ	4.13	70.7	67.4	77.4	71.0	71.6	1895.0 (331.6 / 1563.4)
AWQ	4.13	70.6	68.0	77.2	71.1	71.7	1888.4 (325.7 / 1562.7)
ClusComp	4.13	70.0	68.7	77.6	71.1	71.8	1915.7 (322.1 / 1593.6)
GPTQ	3.13	66.2	65.1	75.6	67.4	68.6	1831.8 (290.1 / 1541.7)
AWQ	3.13	67.7	65.4	74.4	68.0	68.9	1840.3 (298.6 / 1541.7)
ClusComp	2.87	68.7	65.8	74.8	67.7	69.3	1872.6 (331.1 / 1541.5)
GPTQ	2.13	0.0	0.0	0.0	0.0	0.0	0.0 (0.0 / 0.0)
AWQ	2.13	0.0	0.0	0.0	0.0	0.0	0.0 (0.0 / 0.0)
ClusComp	2.14	53.9	53.1	56.7	50.1	53.5	1673.0 (294.6 / 1378.4)

Table 2: Zero-shot multimodal evaluation, with baseline results from Huang et al. (2024c).

vron et al., 2023a,b; Dubey et al., 2024). We mea-400 sure the performance of compressed LLMs on zero-401 shot and language modeling tasks. For zero-shot 402 evaluation, we apply 6 tasks from lm-eval v0.4.4 403 (Gao et al., 2024), i.e. PIQA (Bisk et al., 2020), 404 ARC-e/c (Clark et al., 2018), BoolQ (Clark et al., 405 2019), HellaSwag (Zellers et al., 2019) and Wino-406 Grande (Sakaguchi et al., 2020). For language 407 modeling, we report the perplexity on the whole 408 test set of WikiText2 (Merity et al., 2017) and on 409 256 samples from the validation set of C4 (Raffel 410 et al., 2020b) with a sequence length of 2048 as 411 our baselines. We also apply ClusComp to LLaVA-412 413 Next-8B (Li et al., 2024b), and evaluate it on 5 multimodal tasks from lmms-eval v0.2.3 (Li et al., 414 2024a) to show its broad applicability, i.e. AI2D 415 (Kembhavi et al., 2016), ChartQA (Masry et al., 416 2022), DocVQA (Mathew et al., 2021), MMBench 417 (Liu et al., 2023) and MME (Yin et al., 2023). 418

> **Baselines.** Here we primarily compare Clus-Comp with three categories of baselines: (1) statistic-based methods without neural training, including vanilla RTN, GPTQ (Frantar et al., 2022), AWQ (Lin et al., 2024), and PB-LLM (Yuan et al., 2024);³ (2) gradient-based methods with neural training (such as block-wise distillation), including OmniQuant (Shao et al., 2024), AffineQuant

419

420

421

422

423

424

425

426

(Ma et al., 2024), and SliM-LLM⁺ (Huang et al., 2024b); and (3) quantized codebook-based methods, including QuIP (Chee et al., 2023) and GPTVQ (van Baalen et al., 2024). All baseline results are directly borrowed from the original works or their follow-up works.

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

Settings. We begin by applying K-means clustering to the weight matrices in all linear layers, referring to this method as ClusComp⁻. Next, we use 128 samples from the WikiText2 training set to minimize block-wise error through codebook training only, which we denote as ClusComp. It is important to note that the majority of the aforementioned baselines utilize resources comparable (GPU memory and the number of calibration samples) to those used in ClusComp. Please refer to §B for all experimental details in this section.

Results. The main language modeling results are shown in Figure 1. At the 4-bit level, all methods are comparable. At bit-widths < 4, Clus-Comp consistently outperforms all baselines, with a larger gap for a lower bit-width. Notably, even at the 2-bit level, ClusComp's perplexity remains within a functional range, < 13 on WikiText2. Figure 3 presents the zero-shot evaluation results, where ClusComp again consistently surpasses all baselines across different bit-widths. Furthermore, ClusComp exhibits significantly less sensitivity to

³ClusComp⁻ is also a statistic-based method.

Method	#Bit	Avg. Acc↑	Method	#Bit	Avg. Acc↑
Llama-2-7B	16.00	64.8	Llama-2-13B	16.00	67.8
QuiP#	2.02	57.5	QuiP#	2.01	61.5
AQLM	2.02	56.5	AQLM	1.97	60.6
ClusComp	2.00	56.6	ClusComp	1.99	62.0
ClusComp ⁺	2.00	57.8	ClusComp ⁺	1.99	63.1
Llama-3-8B	16.00	68.6	Llama-3-70B	16.00	75.4
QuiP	2.00	36.8	QuIP	2.00	48.7
PB-LLM	2.00	38.8	PB-LLM	1.70	44.1
DB-LLM	2.00	51.8	BiLLM	1.10	44.2
ClusComp	2.01	53.6	ClusComp	1.14	39.7
ClusComp ⁺	2.01	57.8	ClusComp ⁺	1.14	51.4

Table 3: Zero-shot accuracy on ultra-low-bit LLMs. Please refer to Table A.5 for accuracy on various tasks.

bit-width variations, as indicated by the flatter slope of its accuracy curve.

455

456

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

486

487

488 489

490

491

492

493

We also compress the Llama-3-8B backbone in LLaVA-Next-8B, and report the zero-short performance in Table 2. On average, ClusComp continues to outperform both GPTQ and AWQ, while using a comparable or even lower number of bits. A particularly noteworthy observation occurs at the 2-bit level, where none of the baselines produce correct outputs, whereas ClusComp retains strong performance. In comparison to the 2-bit results for Llama-3-8B in Figure 3 (Right), this suggests that quantizing multimodal models presents unique challenges, warranting further study.

4.2 Push the limit of model compression

We further enhance the performance of 2-bit LLMs and extend the compression boundary to the 1-bit level through efficient recovery training. This is achieved by optimizing the codebook parameters in an end-to-end next-token prediction task.

Baselines. We include BiLLM (Huang et al., 2024a), which performs effectively at the 1-bit compression level. Additionally, three more resource-intensive PTQ baselines are considered: AQLM (Egiazarian et al., 2024), which utilizes a larger number of calibration samples (4-16M tokens); QuIP# (Tseng et al., 2024) and DB-LLM (Chen et al., 2024), both of which employ model-wise distillation and a larger number of calibration samples (24-48M tokens).

Settings. ClusComp employs only 0.3M tokens for its compression. In this experiment, we further finetune the compressed LLM generated by ClusComp through end-to-end training, optimizing the codebook parameters using 16M tokens from a subset of the RedPajama dataset (Computer, 2023). This extended method is referred to as ClusComp⁺.

Results. We report the zero-shot accuracy of ultra-low-bit LLMs in Table 3. On both Llama-2-

Method	Bit	#Trained	MMLU 5-shot ↑	AGIEval 3-shot ↑
Llama-2-7B	16.00	- 100%	45.9	25.7
Full FT	16.00		45.7	27.0
LoRA (r = 128)	16.00	$\begin{array}{r} 4.9\% \\ 100.0\% \\ 100.0\% \\ 0.9\% \\ 1.4\% \\ 1.4\% \end{array}$	45.5	24.7
GaLore	16.00		45.5	24.4
LISA	16.00		46.2	26.1
ClusComp	4.15		47.0	26.5
ClusComp	2.88		45.1	25.6
ClusComp	2.00		30.7	21.8

Table 4: ClusComp performance against efficient finetuning, with baseline results from Pan et al. (2024).

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

7B and 13B, ClusComp already performs comparably to, or surpasses AQLM and QuiP#. With minimal recovery training, ClusComp⁺ consistently outperforms these baselines on average, with the performance gap increasing for larger LLMs, indicating the robust scalability of ClusComp⁺. On Llama-3-8B, ClusComp already exceeds all baselines, and ClusComp⁺ further widens this margin. On Llama-3-70B, ClusComp⁺ achieves remarkable accuracy at the 1-bit level. Furthermore, when comparing the improvements from ClusComp to ClusComp⁺ across different Llama series, a notably larger performance gain is observed on the Llama-3 models, underscoring the effectiveness of ClusComp⁺ on LLMs with more outliers.

4.3 Finetuning quality and efficiency

We can finetune the compressed LLMs on downstream tasks by only training the codebooks.

General-domain finetuning. We finetune the compressed LLMs on a new version of Alpaca, i.e. Alpaca-GPT4 (Peng et al., 2023), and evaluate them on both MMLU and AGIEval (Zhong et al., 2024). Here we mainly compare ClusComp to some memory-efficient finetuning methods that fully finetune the FP16 version, i.e. GaLore (Zhao et al., 2024) and LISA (Pan et al., 2024). As shown in Table 4, finetuning the compressed LLMs at the 4-bit level from ClusComp outperforms all memory-efficient finetuning methods, and rivals full finetuning. In addition, the finetuned LLMs can be used in a low bit, friendly for inference.

The superior finetuning performance can be attributed to three key advantages of ClusComp: (1) ClusComp introduces smaller compression errors; (2) Unlike QLoRA, where compressed LLMs are frozen during finetuning, ClusComp allows for the model to remain unfrozen by training the codebook parameters, enabling further mitigation of compression errors; (3) The low-rank design of



Figure 4: Weight patterns of a cherry-picked layer in Llama-2-7B. Darker red and blue indicate larger and smaller weight values, respectively. Please refer to §A.4 for the visualization setting and the visualization of more layers.



Figure 5: Memory consumption for recovery training or finetuning. 4-bit LLMs are used for ClusComp here.

LoRA limits its expressiveness. In contrast, updating the codebook in ClusComp is analogous to updating a high-rank weight matrix. In addition, the fixed-code design allows uniform training of all centroids, providing greater expressiveness and can even rival full finetuning.

4.4 Discussion

533

534

535

536

537

539

540

541

542

544

545

548

552

Efficiency. Figure 5 illustrates the memory efficiency of ClusComp during training, with a batch size of 1 and a sequence length of 1024. For the 70B LLM, we apply gradient checkpointing (this is not used for the 7B LLM), while omitting any additional memory-saving techniques.

Finetuning the LLM compressed by ClusComp demonstrates memory efficiency in two key ways: (1) The compressed LLM requires less memory for loading onto the GPU compared to the FP16 model; and (2) Only the codebook parameters, which contain a limited number of trainable parameters (< 1%), are updated, as detailed in Table 4. Consequently, the optimizer state size remains small. ClusComp can serve not only as a model compression technique but also as an effective method for both memory-efficient and parameterefficient finetuning. 553

554

555

556

557

558

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

577

578

579

580

581

582

583

585

Weight distribution. In Figure 4, we compare ClusComp's weight distribution to OmniQuant's. Notably, OmniQaunt applies uniform quantization. Overall, ClusComp's weight distribution in different bit levels can better simulate the original weight distribution, explaining ClusComp's effectiveness.

More results. Due to page limit, we recommend readers to (1) Appendix A.5 for more finetuning results with a comparison to QLoRA, LoftQ, QA-LoRA, GPTQ-LoRA and PEQA; (2) Appendix A.6 for comparing ClusComp to SqueezeLLM and AdaDim; (3) Appendix A.7 to check whether the number of clusters n or the cluster dimension ghas a greater impact on the performance; and (4) Appendix A.8 for ClusComp's throughput.

5 Conclusion

The newly introduced model compression technique, ClusComp, operates by (1) independently applying clustering to the weight matrices to produce both the codebook and corresponding codes, (2) reducing compression error through blockwise knowledge distillation, and (3) enhancing model performance via efficient recovery finetuning. Comprehensive experiments demonstrate its effectiveness as a compression method at 1-4 bit levels, while also showcasing its parameter and memory efficiency for finetuning, with a competitive performance with full finetuning.

Limitations

interested readers.

References

- 588 589
- 5
- 592
- 59
- 594

5

598

5

60

606

607

610

611

612

613

614

615

616

617

618

619

622

623

624

626

627

629

630

631

634

638

301

Yoshua Bengio, Nicholas Léonard, and Aaron C. Courville. 2013. Estimating or propagating gradients through stochastic neurons for conditional computation. *CoRR*, abs/1308.3432.

We only evaluate ClusComp on a limited number

of tasks with a total number of 7 models due to

time and resource limitations. We couldn't guaran-

tee its effectiveness on the other tasks and LLMs,

and are still working on including more tasks and

models to show its generalization. In addition, ap-

plying ClusComp to image generation tasks is not explored in this paper. We leave this exploration to

ClusComp only compresses the weight matrices,

while leaving the activations in FP16, which might

cause OOM error for the long-context generation.

However, we can combine ClusComp with activa-

tion quantization methods to address this issue.

- Dan Biderman, Jose Javier Gonzalez Ortiz, Jacob Portes, Mansheej Paul, Philip Greengard, Connor Jennings, Daniel King, Sam Havens, Vitaliy Chiley, Jonathan Frankle, Cody Blakeney, and John P. Cunningham. 2024. Lora learns less and forgets less. *CoRR*, abs/2405.09673.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2020. PIQA: reasoning about physical commonsense in natural language. In The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020, pages 7432– 7439. AAAI Press.
- Yelysei Bondarenko, Markus Nagel, and Tijmen Blankevoort. 2023. Quantizable transformers: Removing outliers by helping attention heads do nothing. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei.

2020. Language models are few-shot learners. *CoRR*, abs/2005.14165.

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

- Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher De Sa. 2023. Quip: 2-bit quantization of large language models with guarantees. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023.
- Hong Chen, Chengtao Lv, Liang Ding, Haotong Qin, Xiabin Zhou, Yifu Ding, Xuebo Liu, Min Zhang, Jinyang Guo, Xianglong Liu, and Dacheng Tao. 2024.
 DB-LLM: accurate dual-binarization for efficient Ilms. In Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024, pages 8719– 8730. Association for Computational Linguistics.
- 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. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers), pages 2924–2936. Association for Computational Linguistics.
- 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. *CoRR*, abs/1803.05457.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *CoRR*, abs/2110.14168.
- Together Computer. 2023. Redpajama: An open source recipe to reproduce llama training dataset.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, and et al. 2024. Deepseek-v3 technical report. *CoRR*, abs/2412.19437.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023.
- Tim Dettmers, Ruslan Svirschevski, Vage Egiazarian, Denis Kuznedelev, Elias Frantar, Saleh Ashkboos, Alexander Borzunov, Torsten Hoefler, and Dan Alistarh. 2024. Spqr: A sparse-quantized representation for near-lossless LLM weight compression. In *The*

807

808

Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024. OpenReview.net.

695

696

701

705

706

707

708

710

711

712

713

714

715

717

718

719

720

721

724

727

728

730

731

732

733

735

736

737

740

741

742

743

744

745

746

747

748

- Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The faiss library. *CoRR*, abs/2401.08281.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and et al. 2024. The llama 3 herd of models. *CoRR*, abs/2407.21783.
- Vage Egiazarian, Andrei Panferov, Denis Kuznedelev, Elias Frantar, Artem Babenko, and Dan Alistarh. 2024. Extreme compression of large language models via additive quantization. In Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024. OpenReview.net.
- Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. 2022. GPTQ: accurate post-training quantization for generative pre-trained transformers. *CoRR*, abs/2210.17323.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2024. A framework for few-shot language model evaluation.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net.
- Jung Hwan Heo, Jeonghoon Kim, Beomseok Kwon, Byeongwook Kim, Se Jung Kwon, and Dongsoo Lee. 2024. Rethinking channel dimensions to isolate outliers for low-bit weight quantization of large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net.
- 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. In Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024. OpenReview.net.

- Wei Huang, Haotong Qin, Yangdong Liu, Yawei Li, Xianglong Liu, Luca Benini, Michele Magno, and Xiaojuan Qi. 2024b. Slim-llm: Salience-driven mixedprecision quantization for large language models. *CoRR*, abs/2405.14917.
- Wei Huang, Xingyu Zheng, Xudong Ma, Haotong Qin, Chengtao Lv, Hong Chen, Jie Luo, Xiaojuan Qi, Xianglong Liu, and Michele Magno. 2024c. An empirical study of llama3 quantization: From llms to mllms. *Preprint*, arXiv:2404.14047.
- Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Min Joon Seo, Hannaneh Hajishirzi, and Ali Farhadi. 2016. A diagram is worth a dozen images. In Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part IV, volume 9908 of Lecture Notes in Computer Science, pages 235–251. Springer.
- Jeonghoon Kim, Jung Hyun Lee, Sungdong Kim, Joonsuk Park, Kang Min Yoo, Se Jung Kwon, and Dongsoo Lee. 2023. Memory-efficient fine-tuning of compressed large language models via sub-4-bit integer quantization. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023.
- Sehoon Kim, Coleman Hooper, Amir Gholami, Zhen Dong, Xiuyu Li, Sheng Shen, Michael W. Mahoney, and Kurt Keutzer. 2024. Squeezellm: Dense-andsparse quantization. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings.
- Diederik P. Kingma and Max Welling. 2014. Autoencoding variational bayes. In 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings.
- Bo Li, Peiyuan Zhang, Kaichen Zhang, Fanyi Pu, Xinrun Du, Yuhao Dong, Haotian Liu, Yuanhan Zhang, Ge Zhang, Chunyuan Li, and Ziwei Liu. 2024a. Lmms-eval: Accelerating the development of large multimoal models.
- Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. 2024b. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *CoRR*, abs/2407.07895.
- Yixiao Li, Yifan Yu, Chen Liang, Nikos Karampatziakis, Pengcheng He, Weizhu Chen, and Tuo Zhao. 2024c. Loftq: Lora-fine-tuning-aware quantization for large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net.

- 810 811
- 812
- 813 814
- 815 816
- 817 818
- 8
- 821 822
- 823 824
- 825 826
- 8
- 8
- 8
- 831 832 833
- 8
- 83
- 83
- 84 84
- 842

- 852 853
- 8

855 856

> 857 858

8

8

86 86

863 864

- Baohao Liao and Christof Monz. 2024a. 3-in-1: 2d rotary adaptation for efficient finetuning, efficient batching and composability. *Preprint*, arXiv:2409.00119.
- Baohao Liao and Christof Monz. 2024b. Apiq: Finetuning of 2-bit quantized large language model. *CoRR*, abs/2402.05147.
- Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. 2024. AWQ: activation-aware weight quantization for ondevice LLM compression and acceleration. In Proceedings of the Seventh Annual Conference on Machine Learning and Systems, MLSys 2024, Santa Clara, CA, USA, May 13-16, 2024. mlsys.org.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. 2023. Mmbench: Is your multi-modal model an all-around player? *CoRR*, abs/2307.06281.
- Stuart P. Lloyd. 1982. Least squares quantization in PCM. *IEEE Trans. Inf. Theory*, 28(2):129–136.
- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net.
- Yuexiao Ma, Huixia Li, Xiawu Zheng, Feng Ling, Xuefeng Xiao, Rui Wang, Shilei Wen, Fei Chao, and Rongrong Ji. 2024. Affinequant: Affine transformation quantization for large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11,* 2024. OpenReview.net.
- Vladimir Malinovskii, Denis Mazur, Ivan Ilin, Denis Kuznedelev, Konstantin Burlachenko, Kai Yi, Dan Alistarh, and Peter Richtárik. 2024. Pv-tuning: Beyond straight-through estimation for extreme LLM compression. *CoRR*, abs/2405.14852.
- Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq R. Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the Association* for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022, pages 2263–2279. Association for Computational Linguistics.
- Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. 2021. Docvqa: A dataset for VQA on document images. In *IEEE Winter Conference on Applications* of Computer Vision, WACV 2021, Waikoloa, HI, USA, January 3-8, 2021, pages 2199–2208. IEEE.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2017. Pointer sentinel mixture models. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. Open-Review.net.

Rui Pan, Xiang Liu, Shizhe Diao, Renjie Pi, Jipeng Zhang, Chi Han, and Tong Zhang. 2024. LISA: layerwise importance sampling for memory-efficient large language model fine-tuning. *CoRR*, abs/2403.17919. 865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

- Yeonhong Park, Jake Hyun, SangLyul Cho, Bonggeun Sim, and Jae W. Lee. 2024. Any-precision LLM: low-cost deployment of multiple, different-sized llms. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27,* 2024. OpenReview.net.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in pytorch. In *NIPS-W*.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020a. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020b. Exploring the limits of transfer learning with a unified text-to-text transformer. J. Mach. Learn. Res., 21:140:1–140:67.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. Winogrande: An adversarial winograd schema challenge at scale. In *The Thirty-Fourth AAAI Conference on Artificial Intelli*gence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020, pages 8732– 8740. AAAI Press.
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilic, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, and et al. 2022. BLOOM: A 176bparameter open-access multilingual language model. *CoRR*, abs/2211.05100.
- Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang, Peng Gao, Yu Qiao, and Ping Luo. 2024. Omniquant: Omnidirectionally calibrated quantization for large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net.
- Mingjie Sun, Zhuang Liu, Anna Bair, and J. Zico Kolter. 2024a. A simple and effective pruning approach for large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net.

1005

1006

1007

1008

1010

1011

1012

978

979

980

Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan. 2024b. Autoregressive model beats diffusion: Llama for scalable image generation. *CoRR*, abs/2406.06525.

921

922

925

930

931

932

933

934

941

942

943

945

947

948

950

953

955

957

960

961

962

963

965

966 967

968

969

970

971

972

973

974

975

976

- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https:// github.com/tatsu-lab/stanford_alpaca.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, and et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *CoRR*, abs/2307.09288.
- Albert Tseng, Jerry Chee, Qingyao Sun, Volodymyr Kuleshov, and Christopher De Sa. 2024. Quip#: Even better LLM quantization with hadamard incoherence and lattice codebooks. In Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024. OpenReview.net.
- Mart van Baalen, Andrey Kuzmin, Markus Nagel, Peter Couperus, Cédric Bastoul, Eric Mahurin, Tijmen Blankevoort, and Paul N. Whatmough. 2024.
 GPTVQ: the blessing of dimensionality for LLM quantization. *CoRR*, abs/2402.15319.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 5998–6008.
- Yuhui Xu, Lingxi Xie, Xiaotao Gu, Xin Chen, Heng Chang, Hengheng Zhang, Zhengsu Chen, Xiaopeng Zhang, and Qi Tian. 2024a. Qa-lora: Quantizationaware low-rank adaptation of large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May* 7-11, 2024. OpenReview.net.
- Yuhui Xu, Lingxi Xie, Xiaotao Gu, Xin Chen, Heng Chang, Hengheng Zhang, Zhengsu Chen, Xiaopeng Zhang, and Qi Tian. 2024b. Qa-lora: Quantizationaware low-rank adaptation of large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May* 7-11, 2024. OpenReview.net.
- Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. 2023. A survey on multimodal large language models. *CoRR*, abs/2306.13549.

- Zhihang Yuan, Yuzhang Shang, and Zhen Dong. 2024. PB-LLM: partially binarized large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May* 7-11, 2024. OpenReview.net.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? In Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers, pages 4791–4800. Association for Computational Linguistics.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona T. Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. OPT: open pre-trained transformer language models. *CoRR*, abs/2205.01068.
- Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong Tian. 2024. Galore: Memory-efficient LLM training by gradient low-rank projection. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.
- Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. 2024. Agieval: A human-centric benchmark for evaluating foundation models. In Findings of the Association for Computational Linguistics: NAACL 2024, Mexico City, Mexico, June 16-21, 2024, pages 2299–2314. Association for Computational Linguistics.



Figure A.1: From Llama-1 to Llama-3, LLMs exhibit increasing challenges for quantization. **Left:** The average kurtosis of weights across various layers in the Llama series, previously shown in Figure 2 (Right). Please refer to Figure B.1 for the kurtosis of different layer types. **Right:** The average standard deviation of the Wanda score across various layers. Please refer to Figure B.2 for the Wanda scores of different layer types. Both metrics indicate that Llama-3 has higher variance in most layers, reflecting the presence of more outliers and thus greater difficulty for quantization.

A More Results and Discussions

1013

1014

1015

1016

1019

1020

1022

1023

1024

1025

1026

1029

1030

1031

1033

1035

1036

1037

1039

1040

1041

1042

1044

A.1 Quantization difficulty in Llama series

Previous works (Lin et al., 2024; Sun et al., 2024a; Heo et al., 2024) suggest that weight patterns can be identified based on activations rather than relying solely on weight magnitudes. We incorporate an additional analysis based on the Wanda score (Sun et al., 2024a) distribution to illustrate the increasing challenges of quantization in the Llama series.

Given the weight matrix $W \in \mathbb{R}^{d_{out} \times d_{in}}$ of a linear layer and the input activations $X \in \mathbb{R}^{NL \times d_{in}}$, where N and L represent the batch and sequence dimensions, the Wanda score $S \in \mathbb{R}^{d_{out} \times d_{in}}$ is computed as $S_{ij} = |W_{ij}| \cdot ||X_j||_2$. A smaller Wanda score within a row of W (on a per-output basis) indicates a less significant weight element.

In Figure A.1 (Right), we present the standard deviation of the Wanda scores across different layers. The results show that Llama-2 exhibits a larger standard deviation compared to Llama-1, while Llama-3 exceeds Llama-2 in this metric. A higher standard deviation reflects a more dispersed Wanda score distribution, indicating that a greater proportion of weight elements are effective and diverse. Consequently, quantization becomes more challenging, as the expanded distribution stretches the quantization grid.

A.2 Code distribution

As shown in Figure A.2, the distribution of codes after clustering is uniform. Training the codebook alone can have a similar effect as full finetuning since all weight values are updated uniformly.



Figure A.2: Histogram of the codes.

A.3 Full results on language modeling

In Table A.1 and Table A.2, we show the full comparison across different LLMs and bit-levels on two language modeling tasks, WikiText2 and C4.

1046

1047

1048

1050

1052

1053

1056

1057

A.4 Weight visualization

To make the weight pattern more obvious, we apply these sequential processing steps: (1) take the absolute weight values; (2) downsample the grids with 8×8 maxpool kernels; (3) calculate the logarithm of these values; (4) normalize the log-values. We offer the visualization of all layers in the first and last blocks of Llama-2-7B in Figure A.3 and A.4.

A.5 More finetuning results

In-domain finetuning. We finetune Llama-2-7B 1058 on the training sets of WikiText2 and GSM8K (Cobbe et al., 2021), and report the perplexity 1060 and accuracy on their respective validation/test set. 1061 ClusComp is compared with two LoRA-based techniques: QLoRA (Dettmers et al., 2023) and LoftQ (Li et al., 2024c). As shown in Table A.6, Clus-1064 Comp consistently achieves superior results with 1065 fewer bits (except at the 4-bit level on WikiText2, 1066 where it performs comparably to LoftQ), even outperforming LoRA-finetuning of the FP16 model 1068

1090

1091 1092

1094 1095

1096 1097

1098 1099

1100 1101

1102

1103 1104

1105 1106

1107

1108

1109

A.8 Throughput

choose $n < 2^{16}$.

Originally, the data type of the codebook C is in 16-1110 bit, which facilitates our following recovery train-1111 ing and finetuning step. However, if the codebook 1112 1113 can be further quantized, we have two additional advantages: (1) The model size can be slightly re-1114 duced (Only slightly since the majority of bits is 1115 allocated to the code q.); (2) The inference can 1116 speed up, since the codebook becomes smaller. In 1117

on GSM8K with a 2.54-bit compressed LLM.

General-domain finetuning. We finetune the

compressed LLMs on Alpaca-GPT3.5 (Taori et al.,

2023) and evaluate them using the MMLU bench-

mark (Hendrycks et al., 2021). ClusComp is com-

pared against baseline methods that merge trained

LoRA modules into the quantized linear layers af-

ter finetuning, i.e. GPTQ-LoRA, QA-LoRA (Xu

As shown in Table A.7, ClusComp consistently

outperforms the baseline methods across different

LLMs and bit-widths, while using fewer bits. The

only exception occurs at the 4-bit level for Llama-1-

7B, where ClusComp underperforms QA-LoRA by

a small margin of 0.3 accuracy. The performance

gap between ClusComp and baselines is enlarged

In this section, we compare ClusComp against

SqueezeLLM (Kim et al., 2024) and AdaDim (Heo

et al., 2024). As presented in Table A.8, Clus-

Comp consistently achieves lower perplexity than

SqueezeLLM, at comparable or even lower bit pre-

cision. Similarly, as shown in Table A.9, Clus-

Comp outperforms AdaDim on both MMLU and

We conducted experiments to determine whether

the number of clusters n or the cluster dimension q

has a greater impact on the performance of quan-

tized LLMs. As shown in Table A.10, increasing n

positively affects performance more than reducing

q. This finding underpins our choice of $n \approx 2^{16}$.

However, while n plays a crucial role in enhancing

performance, selecting $n > 2^{16}$ would necessitate

using 32-bit storage for the code q, substantially

increasing the bits-per-parameter and adversely af-

fecting memory efficiency. Therefore, we always

A.7 Ablation study on the group size and

number of clusters

for lower bits or recent LLM series.

A.6 More baselines

CSR benchmarks.

et al., 2024b) and PEQA (Kim et al., 2023).

sum, we can keep the codebook in 16-bit if we want to do the recovery training or finetuning. If we are only interested in inference, we can further quantize the codebook to a lower bit level.

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1160

As shown in Table A.11, the performance doesn't change if we quantize the codebook from 16-bit to 8-bit. When quantizing the codebook to 4-bit, the perplexity slightly increases, but is still comparable to the best baseline at the 4-bit level and outperforms the best baseline at the 2-bit level. However, if we further quantize the codebook to the 2-bit level, the perplexity increases significantly. Therefore, we can safely quantize the codebook to 8-bit or 4-bit.

In Table A.12, we show the throughput of Clus-Comp. When the codebook stays in the FP16 format, the throughput is only slightly better than the FP16 LLM, because non-uniform compression is not friendly to hardware. To enhance the inference, we can further quantize the codebook to a lower bit level. When the codebook is quantized to a 4-bit level, we can roughly achieve a $2 \times$ speedup.

A.9 **Compression performance of the** multimodal backbone

In Table A.13, we show the performance of the compressed Llama-3-8B backbone in LLaVA-Next-8B. ClusComp again shows its superior performance.

Experimental Details B

B.1 Clustering

We use the K-means implementation from the Faiss library (Douze et al., 2024). The number of iterations is set to 20, with all default settings for other arguments.

B.2 Block-wise error minimization

For all LLMs, 128 calibration sentences with a 1153 length of 2048 tokens are randomly selected from 1154 the WikiText-2 training set (Merity et al., 2017). 1155 The detailed hyper-parameters are listed in Table 1156 B.1. Only the codebooks are trained while keeping 1157 all other parameters (from the embedding layer, 1158 output layer and normalization layers) frozen. 1159

B.3 Recovery training

For the recovery training, we randomly sample 1161 1024 sentences with a length of 4096 tokens from 1162 1163a subset of RedPajama (Computer, 2023)4. Then1164we train the compressed LLMs to predict the next1165token by only tuning the codebook parameters. The1166hyperparameters used in this step are listed in Table1167B.1.

1168 B.4 In-domain finetuning

We follow the settings from (Li et al., 2024c), and 1169 finetune the compressed LLM on the training set 1170 of WikiText2 and on the training set of GSM8K. 1171 The hyperparameters for finetuning are listed in 1172 Table B.2. We evaluate the finetuned model on the 1173 validation set of WikiText2 and on the test set of 1174 GSM8K every epoch and report the best perplexity 1175 or accuracy. 1176

1177 B.5 General-domain finetuning

1178The finetuning hyper-parameters are listed in Table1179B.2, which is similar to the ones in QA-LoRA (Xu1180et al., 2024b) on Alpaca-GPT3.5, or to the ones in1181LISA (Pan et al., 2024) on Alpaca-GPT4.

⁴https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T-Sample

Method	Setting	#Bit	1-7B	2-7B	2-13B	2-70B	3-8B	3-70B
-	-	16.00	5.68	5.47	4.88	3.31	6.12	2.90
RTN	w4g128	4.13	5.96	5.72	4.98	3.46	8.50	3.60
GPTQ	w4g128	4.13	5.85	5.61	4.98	3.42	6.50	3.30
AWQ	w4g128	4.13	5.81	5.62	4.97	-	6.60	3.30
GPTVQ	w4g128	4.13	-	5.68	4.97	3.39	-	-
AffineQuant	w4g128	4.13	5.77	5.58	4.95	-	-	-
OmniQuant	w4g128	4.16	5.77	5.58	4.95	3.40	-	-
ClusComp [—] ClusComp	g4n65500 g4n65500	≤ 4.14 ≤ 4.14	5.88 (4.14) 5.73 (4.14)	5.67 (4.14) 5.54 (4.14)	5.04 (4.09) 4.94 (4.09)	3.44 (4.03) 3.40 (4.03)	6.59 (4.13) 6.39 (4.13)	3.28 (4.03) 3.12 (4.03)
RTN	w4	4.00	6.43	6.11	5.20	3.67	8.70	-
GPTQ	w4	4.00	6.13	5.83	5.13	3.58	7.00	-
AWQ	w4	4.00	6.08	6.15	5.12	-	7.10	-
QuIP	w4	4.00	-	-	-	3.53	-	-
AffineQuant	w4	4.00	5.84	5.69	5.01	-	-	-
OmniQuant	W4	4.00	5.86	5.74	5.02	3.47	-	-
ClusComp	g5n65500	3.38	5.27 5.84	5.90 5.67	_	_	-	-
	551105500	0.00	2.04	2.07				
RTN	w3g128	3.13	7.01	6.66	5.51	3.97	27.91	11.84
GPTQ	w3g128	3.13	6.55	6.29	5.42	3.85	8.22	5.22
AWQ	w3g128	3.13	6.46	6.24	5.32	-	8.19	4.81
SliM-LLM	w3g128	3.13	6.07	5.94	5.11	3.35	-	-
GPTVO	w3g128	3.13	0.14	0.08 5 82	5.28	- 2 55	-	-
OmniQuant	w3g128 w3g128	3.15	- 6 15	5.62 6.03	5.10	3.78	-	-
	#5g120	5.15	0.15	0.05	5.20	5.76		
RTN	w3	3.00	25.73	5.4e2	10.68	7.52	2.2e3	-
GPTQ	w3	3.00	8.06	8.37	6.44	4.82	13.0	-
AWQ	w3	3.00	11.88	24.00	10.45	-	12.8	-
AffinaQuant	W3	3.00	- 6 20	- 6 55	- 5.62	3.85	1.5	-
OmniQuant	w3	3.00	6.49	6.58	5.58	- 3.92	-	-
ClusComp ⁻	g6n65500	< 2.89	6 74 (2 89)	6 54 (2 89)	6 27 (2 81)	4 02 (2 72)	8 77 (2 87)	4 98 (2 72)
ClusComp	g6n65500	≤2.89	6.01 (2.89)	5.86 (2.89)	5.18 (2.81)	3.72 (2.72)	7.34 (2.87)	4.63 (2.72)
ClusComp ⁻	o7n65500	2.54	7 79	7 64	_	_	-	-
ClusComp	g7n65500	2.54	6.28	6.15	-	-	-	-
-	2 ()	0.05	10.0	4.2.2	26.22	10.21		
RIN	w2g64	2.25	1.9e2	4.3e2	26.22	10.31 NAN	-	-
AWO	w2g04 w2g64	2.23	22.10	20.85	1 205	INAIN	1.862	-
GPTVO	w2g04 w2g64	2.25	-	7 22	6.08	4 39	-	-
AffineOuant	w2g64	2.25	8.35	9.05	7.11	-	-	-
OmniQuant	w2g64	2.28	8.90	9.62	7.56	6.11	-	-
ClusComp ⁻	g8n65500	≤ 2.29	10.25 (2.29)	11.10 (2.29)	14.39 (2.19)	-	29.20 (2.27)	-
ClusComp	g8n65500	≤2.29	6.66 (2.29)	6.61 (2.29)	5.74 (2.19)	-	9.68 (2.27)	-
RTN	w2g128	2.13	1.9e3	4.2e3	1.2e2	27.27	1.9e3	4.6e5
GPTQ	w2g128	2.13	44.01	36.77	28.14	NAN	2.1e2	11.9
AWQ	w2g128	2.13	2.665	2.265	1.2e5	-	1./e6	1./e6
SIIWI-LLM	w2g128	2.13	9.08	10.87	1.59	0.44 6.64	- 84.07	-
PB-LIM	w2g128	2.13	-	25 37	49.81	0.04 NAN	64.97 44.12	11.68
GPTVO	w2g128	2.13	_	8 23	6 50	4 64	-	-
AffineQuant	w2g128	2.13	13.51	10.87	7.64	-	-	
OmniQuant	w2g128	2.14	9.72	11.06	8.26	6.55	-	-
ClusComp ⁻	g8	≤ 2.15	28.76	21.9	14.50	5.43	2.1e2	11.40
ClusComp	g8	≤2.15	7.06 (2.15,n35000)	7.04 (2.15,n35000)	5.85 (2.14,n50000)	4.37 (2.07,n65500)	11.57 (2.14,n35000)	7.61 (2.07,n65500)
ClusCome -	c0n65500	2.11		22.71				
ClusComp	g9n65500 g9n65500	2.11 2.11	_	7.12	-	-	-	-
DTN		2.00	1.1.5	2.9-4	5 (-1	2.0-4	27-6	
GPTO	W2 W2	2.00	1.1e5 2.1o2	5.8e4	5.6e4 2.1o2	2.0e4 77.05	2.7eb	-
OnIP	w∠ w2	2.00	2.103	1.105	2.103	633	3.704 85.1	-
AffineOuant	w2	2.00	9.53	35.07	12.42	-	-	-
OmniQuant	w2	2.00	15.47	37.37	17.21	7.81	-	-
ClusComp ⁻	g9	≤2.01	65.09	52.38	22.90	9.84	3.1e2	-
ClusComp	g9	≤ 2.01	7.49	7.50	6.17	4.83	12.33	-
			(2.00,n45000)	(2.00,n45000)	(1.99,n65500)	(1.85,n65500)	(2.01,n50000)	-

Table A.1: The full perplexity results of Llama series on WikiText2. "g" and "n" denote the dimension and number of centroids in the codebook, respectively. The number in the brackets is the exact bits of different settings for different LLMs.

Method	Setting	#Bit	1-7B	2-7B	2-13B	2-70B	3-8B	3-70B
-	-	16.00	7.08	6.97	6.46	5.52	9.20	5.87
RTN GPTQ AWQ AffineQuant OmniQuant	w4g128 w4g128 w4g128 w4g128 w4g128 w4g128	4.13 4.13 4.13 4.13 4.13 4.16	7.37 7.21 7.21 7.20 7.21	7.24 7.12 7.13 7.12 7.12	6.58 6.56 6.56 6.56 6.56	5.63 5.58 - - 5.58	13.40 10.40 9.40	8.90 6.94 7.00
ClusComp ⁻ ClusComp	g4n65500 g4n65500	≤4.14 ≤4.13	7.27 (4.14) 7.17 (4.14)	7.16 (4.14) 7.09 (4.14)	6.63 (4.09) 6.55 (4.09)	5.61 (4.03) 5.61 (4.03)	9.39 (4.13) 9.27 (4.13)	7.02 (4.03) 6.99 (4.03)
RTN GPTQ AWQ QuIP AffineQuant OmniQuant ClusComp ClusComp	w3 w3 w3 w3 w3 w3 g6n65500 g6n65500	3.00 3.00 3.00 3.00 3.00 3.00 ≤2.89 ≤ 2.89	28.26 9.49 13.26 - 8.03 8.19 8.14 (2.89) 7.64 (2.89)	4.0e2 9.81 23.85 - 8.57 8.65 8.19 (2.89) 7.61 (2.89)	12.51 8.02 13.07 - 7.56 7.44 8.21 (2.81) 6.91 (2.81)	10.02 6.57 - 6.14 - 6.06 6.06 (2.72) 5.86 (2.72)	2.2e3 13.0 12.8 - - 12.41 (2.87) 11.31 (2.87)	8.26 (2.72) 8.26 (2.72)
ClusComp ⁻ ClusComp	g7n65500 g7n65500	2.54 2.54	9.46 8.10	9.51 8.13	-	-	-	-
RTN GPTQ AWQ OmniQuant ClusComp	w2g64 w2g64 w2g64 w2g64 g8n65500	$2.25 \\ 2.25 \\ 2.25 \\ 2.28 \\ \leq 2.29 \\ $	1.5e2 17.71 2.8e5 11.78 13.06 (2.29)	4.8e2 19.40 1.6e5 12.72 14.07 (2.29)	28.69 12.48 9.5e4 10.05 19.75 (2.19)	13.43 NAN - 7.88	38.68 (2.27)	- - - -
RTN GPTQ AWQ SliM-LLM ⁺ QuIP PB-LLM AffineQuant OmniOuant	w2g128 w2g128 w2g128 w2g128 w2g128 w2g128 w2g128 w2g128 w2g128 w2g128	2.13 2.13 2.13 2.13 2.13 2.13 2.13 2.13	8.76 (2.29) 1.0e3 27.71 1.9e5 14.99 - - 12.97	4.9e3 33.70 1.7e5 18.18 31.94 29.84 16.02 15.02	1.4e2 20.97 9.4e4 10.24 16.16 19.82 10.98 11.05	- 42.13 NAN - 8.40 8.17 8.95 - 8.52	1.9e3 2.1e2 1.7e6 - 1.3e2 79.21	- - - 22.24 33.91
ClusComp ⁻ ClusComp	g8 g8	≤2.15 ≤2.15	29.67 9.33 (2.15,n35000)	25.26 9.49 (2.15,n35000)	18.83 7.92 (2.14,n50000)	7.59 6.44 (2.07,n65500)	1.9e2 17.89 (2.14,n35000)	16.52 10.81 (2.07,n65500)
ClusComp ⁻ ClusComp	g9n65500 g9n65500	2.11 2.11	-	27.37 9.73	-	-	-	-
RTN GPTQ QuIP OmniQuant ClusComp ⁻ ClusComp	w2 w2 w2 g9 g9	$2.00 \\ 2.00 \\ 2.00 \\ 2.00 \\ \le 2.01 \\ \le 2.01$	1.3e5 6.9e2 - 24.89 74.61 10.11	4.8e4 NAN - 90.64 50.08 10.29	7.2e4 3.2e2 - 26.76 24.47 8.49	2.4e4 48.82 - 12.28 13.96 7.02	2.7e6 5.7e4 1.3e2 8.2e5 2.2e2 21.45	- - -
			(2.00,n45000)	(2.00,n45000)	(1.99,n65500)	(1.85,n65500)	(2.01,n50000)	-

Table A.2: The full perplexity results of Llama series on C4. "g" and "n" denote the dimension and number of centroids in the codebook, respectively. The number in the brackets is the exact bits of different settings for different LLMs.

		P	ÍQA	Al	RC-e	AI	RC-c	BoolQ	Hell	aSwag	WinoGrande	
Method	#Bit	acc	acc_n	acc	acc_n	acc	acc_n	acc	acc	acc_n	acc	Avg
Llama-2-7B	16.00	-	79.1	-	74.6	-	46.3	77.7	-	76.0	69.1	70.5
ClusComp	4.14	77.5	79.2	75.3	72.8	42.7	45.3	76.2	56.5	75.1	68.9	69.6
RTN	3.13	-	76.8	-	70.5	-	42.9	71.7	-	74.0	67.6	67.3
GPTQ	3.13	-	77.4	-	68.1	-	40.7	71.0	-	72.5	67.3	66.2
GPTVQ	3.13	-	77.6	-	72.7	-	43.7	71.7		72.7	67.6	67.7
ClusComp	2.89	76.8	77.6	74.4	71.3	42.3	42.9	74.6	54.4	72.4	68.8	67.9
RTN	2.25	-	58.8	-	36.7	-	24.8	41.9	-	40.4	51.9	42.4
GPTQ	2.25	-	60.8	-	39.0	-	25.2	59.3	-	45.8	55.5	47.6
GPTVQ	2.25	-	73.3	-	63.4	-	35.9	66.3	-	63.9	66.1	61.5
ClusComp	2.29	74.9	76.0	69.8	65.2	37.7	37.4	73.0	51.1	68.4	65.0	64.1
ClusComp	2.15	74.3	75.1	69.6	65.2	35.7	38.4	69.5	49.2	66.4	63.5	63.0
RTN	2.13	-	51.1	-	28.0	-	25.0	41.1	-	26.6	49.9	36.9
GPTQ	2.13	-	54.8	-	30.6	-	25.1	53.4	-	33.1	51.5	41.4
GPTVQ	2.13	-	70.7	-	58.1	-	31.5	63.7	-	58.5	60.9	57.2
ClusComp	2.00	72.6	73.7	67.0	62.8	32.9	36.6	70.9	47.2	63.5	63.4	61.8
Llama-2-13B	16.00	-	80.5	-	77.5	-	49.2	80.5	-	79.4	72.1	73.2
ClusComp	4.09	78.9	79.9	78.9	76.9	47.7	49.2	81.4	60.0	79.0	72.4	73.1
RTN	3.13	-	78.9	-	74.3	-	46.8	77.3	-	76.5	70.8	70.8
GPTQ	3.13	-	79.3	-	75.8	-	47.0	78.9	-	77.2	70.4	71.4
GPTVQ	3.13	-	79.4	-	75.3	-	48.1	79.0	-	77.0	71.7	71.8
ClusComp	2.81	78.7	79.7	78.5	76.7	45.9	47.8	80.7	58.3	76.8	71.4	72.2
RTN	2.25	-	61.6	-	41.6	-	25.4	49.8	-	48.2	51.9	46.4
GPTO	2.25	-	70.1	-	56.7	-	31.6	51.1	-	56.6	58.9	54.2
GPTVO	2.25	-	76.2	-	71.9	-	43.3	67.6	-	70.0	68.2	66.2
ClusComp	2.19	76.6	77.3	75.0	72.9	40.8	43.9	78.1	55.3	73.3	68.4	69.0
ClusComp	2.14	76.7	77.1	73.5	71.6	39.9	42.8	77.5	54.6	73.1	68.0	68.4
RTN	2.13	-	58.4	-	32.3	-	25.5	47.9	-	39.4	48.9	42.1
GPTQ	2.13	-	59.5	-	40.2	-	27.7	57.1	-	41.6	53.4	46.6
GPTVO	2.13	-	75.2	-	68.3	-	39.5	70.7	-	65.7	67.5	64.5
ClusComp	1.99	75.6	77.7	74.7	73.6	39.9	42.1	74.0	53.0	71.0	67.1	67.6

Table A.3: Zero-shot evaluation of the quantized Llama-2-7B and Llama-2-13B, with baseline results taken from van Baalen et al. (2024). "acc" and "acc_n" mean accuracy and normalized accuracy, respectively. We offer the results of all metrics for a convenient comparison of the follow-up works. But only the highlighted **metrics** are used to calculate the average accuracy.

		P	IQA	AI	RC-e	AI	RC-c	BoolQ	Hell	aSwag	WinoGrande	
Method	#Bit	acc	acc_n	acc	acc_n	acc	acc_n	acc	acc	acc_n	acc	Avg
Llama-3-8B	16.00	79.9	-	80.1	-	50.4	-	-	60.2	-	72.8	68.6
RTN	4.13	76.6	-	70.1	-	45.0	-	-	56.8	-	71.0	63.9
GPTQ	4.13	78.4	-	78.8	-	47.7	-	-	59.0	-	72.6	67.3
AWQ	4.13	79.1	-	79.7	-	49.3	-	-	59.1	-	74.0	68.2
SliM-LLM	4.13	78.9	-	79.9	-	49.4	-	-	58.7	-	72.6	67.9
ClusComp	4.13	79.1	80.5	80.9	79.6	49.7	54.1	81.1	59.3	78.3	72.9	68.4
RTN	3.13	62.3	-	32.1	-	22.5	-	-	29.1	-	54.7	40.2
GPTQ	3.13	74.9	-	70.5	-	37.7	-	-	54.3	-	71.1	61.7
AWQ	3.13	77.7	-	74.0	-	43.2	-	-	55.1	-	72.1	64.4
SliM-LLM	3.13	77.8	-	73.7	-	42.9	-	-	55.5	-	72.8	64.5
RTN	3.00	56.2	-	31.1	-	20.0	-	-	27.5	-	53.1	35.6
GPTQ	3.00	60.8	-	38.8	-	22.3	-	-	41.8	-	60.9	44.9
AWQ	3.00	71.9	-	66.7	-	35.1	-	-	50.7	-	64.7	57.8
QuIP	3.00	76.8	-	72.9	-	41.0	-	-	55.4	-	72.5	63.7
ClusComp	2.87	77.7	78.8	76.0	74.5	43.9	47.6	79.0	56.0	74.6	71.0	64.9
ClusComp	2.27	70.6	71.8	63.5	57.4	31.4	35.5	74.7	49.6	66.1	67.1	56.4
RTN	2.13	53.1	-	24.8	-	22.1	-	-	26.9	-	53.1	36.0
GPTQ	2.13	53.9	-	28.8	-	19.9	-	-	27.7	-	50.5	36.2
AWQ	2.13	52.4	-	24.2	-	21.5	-	-	25.6	-	50.7	34.9
SliM-LLM	2.13	57.1	-	35.4	-	26.1	-	-	28.9	-	56.6	40.8
PB-LLM	2.13	57.0	-	37.8	-	17.2	-	-	29.8	-	52.5	38.8
ClusComp	2.14	68.0	67.1	54.7	49.0	26.4	28.8	71.5	47.0	63.0	62.4	51.7
RTN	2.00	53.1	-	24.7	-	21.9	-	-	25.6	-	51.1	35.3
GPTQ	2.00	52.8	-	25.0	-	20.5	-	-	26.6	-	49.6	34.9
AWQ	2.00	55.2	-	25.2	-	21.3	-	-	25.4	-	50.4	35.5
QuIP	2.00	52.9	-	29.0	-	21.3	-	-	29.2	-	51.7	36.8
ClusComp	2.01	70.1	69.6	63.3	57.7	31.9	34.2	66.6	44.4	58.0	58.4	53.6

Table A.4: Zero-shot evaluation of the quantized Llama-3-8B, with baseline results taken from (Huang et al., 2024c). "acc" and "acc_n" mean accuracy and normalized accuracy, respectively. We offer the results of all metrics for a convenient comparison of the follow-up works. But only the highlighted **metrics** (excluding BoolQ) are used to calculate the average accuracy.

Method	#Bit	PIQA	ArcE	ArcC	Hella.	Wino.	Avg
Llama-2-7B	16.00	78.1	76.3	43.4	57.1	69.1	64.8
QuIP#	2.02	75.1	64.6	34.6	48.3	64.9	57.5
AQLM	2.02	73.6	61.9	33.3	49.5	64.2	56.5
ClusComp	2.00	72.6	67.0	32.9	47.2	63.4	56.6
ClusComp ⁺	2.00	73.5	67.2	33.9	49.3	65.1	57.8
Llama-2-13B	16.00	79.1	79.4	48.4	60.0	72.2	67.8
QuIP#	2.01	77.3	69.3	39.5	53.4	67.7	61.5
AQLM	1.97	76.2	69.8	37.8	53.7	65.4	60.6
ClusComp	1.99	75.6	74.7	39.9	53.0	67.1	62.0
ClusComp ⁺	1.99	76.8	74.9	40.7	54.5	68.4	63.1
Llama-3-8B	16.00	79.7	80.1	50.4	60.2	72.6	68.6
QuiP	2.00	52.9	29.0	21.3	29.2	51.7	36.8
PB-LLM	2.00	57.0	37.8	17.2	29.8	52.5	38.8
DB-LLM	2.00	68.9	59.1	28.2	42.1	60.4	51.8
ClusComp	2.01	70.1	63.3	31.9	44.4	58.4	53.6
ClusComp ⁺	2.01	74.8	66.7	34.8	49.6	63.4	57.8
Llama-3-70B	16.00	82.5	86.7	60.4	66.3	80.9	75.4
QuIP	2.00	65.3	48.9	26.5	40.9	61.7	48.7
PB-LLM	1.70	56.5	49.9	25.8	34.9	53.1	44.1
BiLLM	1.10	58.2	46.4	25.1	37.5	53.6	44.2
ClusComp	1.14	56.9	32.5	20.6	32.3	51.6	39.7
ClusComp ⁺	1.14	69.5	57.2	30.1	44.2	56.0	51.4

Table A.5: Zero-shot accuracy on ultra-low-bit LLMs.

Method	#Bit	WikiText2 PPL↓	GSM8K Acc ↑
LoRA	16.00	5.08	36.9
QLoRA LoftQ ClusComp	4.25 + 0.40 4.25 + 0.40 4.15	5.70 5.24 5.26	35.1 35.0 41.0
QLoRA LoftQ	3.25 + 0.40 3.25 + 0.40	5.73 5.63	32.1 32.9
ClusComp	3.38	5.37	39.9
QLoRA LoftQ	2.25 + 0.40 2.25 + 0.40 2.54	NAN 7.85	NAN 20.9
ClusComp	2.34 2.29	6.10	36.0

Table A.6: In-domain finetuning performance on Llama-2-7B. Two bits are shown for baselines since the LoRA modules aren't merged to the quantized LLMs. The first and second numbers denote the quantized LLM and the converted bits from LoRA modules.

			MMLU	(0-shot, A	Acc ↑)			MMLU	(5-shot, A	Acc ↑)	
Method	#Bit	Hums.	STEM	Social	Other	Avg	Hums.	STEM	Social	Other	Avg
Llama-1-7B	16.00	32.4	26.6	31.4	37.2	32.1	33.3	29.8	37.8	38.0	34.6
GPTQ-LoRA QA-LoRA	4.50 4.50	35.7 36.9	30.9 31.4	38.0 40.3	44.0 44.9	37.1 38.3	33.8 36.6	31.3 32.4	37.4 44.8	42.2 44.9	36.0 39.4
PEQA	4.00	-	-	-	-	-	34.9	28.9	37.5	40.1	34.8
ClusComp	4.15	36.9	31.4	40.0	44.2	38.0	36.8	34.1	42.7	43.9	39.1
GPTQ-LoRA	3.50	31.5	28.9 34.1	31.8	36.8	32.2	31.6	30.1	35.6 41 5	39.8 42.7	34.0 37.4
ClusComp	3.38	38.2	32.7	41.2	4 2.3 45.4	39.2	36.3	31.4	41.3	43.0	37.4
GPTQ-LoRA	2.50	24.1	22.1	22.5	23.7	23.2	23.4	26.2	26.4	28.4	25.8
QA-LoRA	2.50	26.4	25.5	25.6	28.7	26.5	27.3	26.1	26.1	30.3	27.5
ClusComp	2.29	32.6	29.7	34.4	37.0	33.3	31.1	30.1	37.8	37.2	33.7
Llama-2-7B	16.00	38.9	32.9	46.6	44.9	40.7	43.0	36.4	51.4	52.2	45.5
QA-LoRA	4.50	41.1	35.4	50.2	50.1	43.9	42.1	34.4	49.1	50.3	43.9
ClusComp	4.15	41.6	36.3	52.3	51.1	44.9	42.8	38.1	52.2	53.1	46.1
Llama-2-13B	16.00	48.1	42.7	60.5	59.5	52.3	53.3	44.1	63.3	61.0	55.3
QA-LoRA	4.50	48.2	41.7	60.4	58.7	51.9	48.0	43.0	59.7	57.4	51.7
ClusComp	4.09	49.2	42.9	61.6	60.2	52.9	52.4	43.2	62.9	61.6	54.7

Table A.7: General-domain finetuning performance, with baseline results from Xu et al. (2024b).

Method	#Bit	Llama-2-7B	Llama-2-13B	Llama-2-70B
-	16.00	5.47	4.88	3.31
SqueezeLLM	4.27	5.57	4.96	-
ClusComp	\leq 4.14	5.54 (4.14)	4.94 (4.09)	-
SqueezeLLM	3.02	6.18	5.36	3.77
ClusComp	\leq 2.89	5.86 (2.89)	5.18 (2.81)	3.72 (2.72)
SqueezeLLM	2.22	10.79	7.91	4.99
SqueezeLLM	2.05	13.64	8.56	5.38
SqueezeLLM	2.01	35.49	41.02	9.44
ClusComp	≤ 2.00	7.50 (2.00)	6.17 (1.99)	4.83 (1.85)

Table A.8: The perplexity of Llama-2 on WikiText2. The values in the brackets are the exact bits of ClusComp for different LLMs. The SqueezeLLM results are taken from Kim et al. (2024).

		Llama-	2-7B	Llama-2-13B		
Method	#Bit	MMLU (5-shot ↑)	\mathbf{CSR} (0-shot \uparrow)	MMLU (5-shot ↑)	$CSR (0-shot \uparrow)$	
-	16.00	46.0	67.9	55.6	70.3	
GPTQ-AdaDim ClusComp	4.13 ≤ 4.14	45.3 45.6	67.7 68.2	54.6 55.1	70.1 70.8	
GPTQ-AdaDim ClusComp	3.13 ≤ 2.89	41.3 43.2	66.4 67.2	52.3 52.3	68.7 69.5	

Table A.9: The accuracy of quantized LLMs on MMLU and four commonsense reasoning (CSR) tasks (PIQA, HellaSwag, WinoGrande and ARC-easy). Following AdaDim, we use Im-eval v0.3.0 (Gao et al., 2024) for the evaluation. The GPTQ-AdaDim results are taken from Heo et al. (2024).

Setting	#Bit	Wiki2
g4	4.14	5.67
g5	3.38	5.90
g6	2.89	6.54
g7	2.54	7.64
g8	2.29	11.10

g7n16384 g7n4096	2.35 2.30	23.14 9.4e2
g8n65500	2.29	11.10
g8n50000	2.15	21.90
g8n4096	2.02	5.4e3

#Bit

Wiki2

Setting

Setting	#Bit	Wiki2
g7n16384	2.35	23.14
g7n4096	2.30	9.4e2
g8n65500	2.29	11.10

(a) n = 65500. Perplexity changes smoothly.

(b) Same *g*. Perplexity changes dramatically.

(c) Same *g*. Perplexity changes dramatically.

Table A.10: Ablation study of ClusComp⁻ on the number of clusters n and the cluster dimension g in the codebook reveals that n plays a more significant role in the performance of the quantized LLM. (a) The perplexity remains relatively stable with variations in g, although changes in g lead to substantial differences in the bit requirement, as most bits are used to store the codes q. Specifically, smaller g values result in larger q. Refer to Table 1 for detailed examples. (b) In contrast, perplexity is highly sensitive to changes in n. Adjusting n causes only a minor change in the bit requirement, as storing the codebook is memory-efficient. (c) For comparable bit budgets, n has a greater impact on performance than q.

Method	Bit for codebook	Avg. Bit	2-7B	2-70B	3-8B	3-70B
-	-	16.00	5.47	3.31	6.12	2.90
Best baseline ClusComp ClusComp	16 8	$\begin{array}{c} 4.13\\ \leq 4.14\\ \leq 4.11\end{array}$	5.58 5.54 (4.14) 5.54 (4.11)	3.39 3.40 (4.03) 3.40 (4.03)	6.50 6.39 (4.13) 6.39 (4.10)	3.30 3.12 (4.03) 3.13 (4.03)
ClusComp ClusComp	4 2	≤ 4.07 ≤ 4.05	5.59 (4.07) 25.44 (4.05)	3.43 (4.02) 5.64 (4.01)	6.52 (4.07) 1.2e5 (4.05)	3.26 (4.02) 96.52 (4.01)
Best baseline ClusComp ClusComp ClusComp ClusComp	16 8 4 2	$ \leq 2.13 \\ \leq 2.07 \\ \leq 2.04 \\ \leq 2.03 \\ \leq 2.02 $	35.07 (2.00) 7.50 (2.00) 7.50 (1.92) 7.63 (1.86) 6.5e3 (1.83)	4.64 (2.13) 4.37 (2.07) 4.37 (2.04) 4.42 (2.03) 21.07 (2.02)	85.10 (2.00) 12.33 (2.01) 12.33 (1.92) 12.77 (1.86) 1.7e5 (1.83)	11.68 (2.13) 7.61 (2.07) 7.63 (2.04) 7.64 (2.03) 2.3e4 (2.02)

Table A.11: The perplexity of ClusComp with quantized codebook on WikiText2. The results of the best baseline are taken from Table A.1. The values in the brackets are the exact bits for different LLMs. We can observe: (1) 8-bit codebook offers the same perplexity as 16-bit's; (2) 4-bit codebook slightly hurts the performance, but is still comparable to the best baseline at the 4-bit level and outperforms the best baseline at the 2-bit level. (3) The results of the 2-bit codebook are not acceptable.

#Bits for Codebook	Avg. #Bit	Llama-2-7B	Llama-2-70B
-	16.00	$1.00 \times$	$1.00 \times$
16.00	2.01	$1.29 \times$	$1.17 \times$
8.00	1.92	$1.31 \times$	$1.19 \times$
4.00	1.86	$2.10 \times$	$1.81 \times$

Table A.12: Throughput comparison between ClusComp and the FP16 version. Originally, the codebook in ClusComp remained in FP16. For further inference speedup, we can uniformly quantize it to a lower bit level.

Method	Setting	Bit	WikiText2↓	$C4\downarrow$	$\textbf{PTB}\downarrow$
-	-	16.00	9.5	14.8	16.3
GPTQ	w4g128	4.13	9.5	14.8	17.1
AWQ	w4g128	4.13	9.9	15.3	16.9
ClusComp ⁻	s4n65500	4.13	9.9	13.6	17.6
ClusComp	s4n65500	4.13	9.7	13.6	17.7
GPTQ	w3g128	3.13	13.0	19.5	28.4
AWQ	w3g128	3.13	11.7	17.9	20.2
ClusComp ⁻	s6n65500	2.87	14.3	16.2	31.9
ClusComp	s6n65500	2.87	10.7	15.3	22.0
GPTQ	w2g128	2.13	83.7	3.1e3	2.0e2
AWQ	w2g128	2.13	1.6e6	2.0e6	2.2e6
ClusComp ⁻	s8n35000	2.14	7.7e2	6.1e3	9.2e2
ClusComp	s8n35000	2.14	14.6	21.8	27.5

Table A.13: The perplexity of the Llama-3-8B backbone in LLaVA-Next-8B, with baseline results from Huang et al. (2024c).

Hyper-parameter	Block-wise error minimization	Recovery training	
Optimizer	AdamW (Loshchilov and Hutter, 20	19; Kingma and Ba, 2015)	
Weight decay	$\{0, 0.1, 0.01\}$	0	
LR	{1e-5, 5e-5, 1e-4, 5e-4}	1e-5	
LR scheduler	constant	cosine	
Warmup ratio	0	0	
Max grad norm	-	0.3	
Sequence length	2048	4096	
Number of samples	128	8192	
Epochs	20	1	
Batch size	8	8	

Table B.1: Hyper-parameters used for the block-wise error minimization and recovery training steps. The underlined settings generally perform well for different scales of LLMs.

Hyper-parameter	WikiText-2	GSM8K	Alpaca-GPT3.5	Alpaca-GPT4
Optimizer	AdamW		AdamW	
Weight decay	0.1		0	
LR	$\{0.7, 1, 3\} \times 10^{-4}$		$\{2, 4, \underline{6}, 8\} \times 10^{-5}$	
LR scheduler	cosine		cosine	
Warmup ratio	3%		6%	
Epochs or max steps	3 epochs	6 epochs	10K steps	2 epochs
Batch size	64	16	- 16	5
Max sequence length 1024 51		512	2048	

Table B.2: Hyperparameters for the finetuning on Llama-2-7B. The underlined settings generally performs well for different bit levels. We report the average performance from three random runs in this paper.



(g) Down projection layer.

Figure A.3: Weight patterns of all layers in the first block (0-th block) of Llama-2-7B. Darker red and blue indicate larger and smaller weight values, respectively. ClusComp's weight distribution of different bit levels can better simulate the original weight distribution.



(g) Down projection layer.

Figure A.4: Weight patterns of all layers in the last block (31-st block) of Llama-2-7B. Darker red and blue indicate larger and smaller weight values, respectively. ClusComp's weight distribution of different bit levels can better simulate the original weight distribution.



Figure B.1: The kurtosis across various layers in different Llama series reveals three key observations: (1) Layers at either the beginning or the end of LLMs tend to exhibit higher kurtosis values; (2) In the majority of layers, the kurtosis follows a consistent trend across Llama series, with Llama-3 showing the highest values, followed by Llama-2, and then Llama-1; (3) Different types of layers display varying scales of kurtosis, suggesting that a bit allocation strategy that accounts for quantization difficulty could yield better results. We leave the exploration of this idea to future work.



Figure B.2: The standard deviation of Wanda score across various layers in different Llama series reveals three key observations: (1) Deeper layers tend to exhibit higher standard deviation; (2) Three layers (query, key and gate) show a clear trend across Llama series, with Llama-3 showing the highest standard deviation, followed by Llama-2, and then Llama-1; (3) The other four layers show a similar standard deviation for all Llama series.

```
class ClusCompLinear(nn.Module):
    def __init__(Self, in_features, out_features, num_clusters, cluster_dim, bias):
        super()._init__()
        self.out_features = out_features
        self.deficiency = out_features % cluster_dim # If the out_features is not dividable by cluster_dim
        if self.deficiency > 0:
            self.deficiency > 0:
            self.deficiency = cluster_dim - self.deficiency) // cluster_dim
        self.codebook = nn.Parameter(torch.empty((num_clusters, cluster_dim), dtype=torch.bfloat16) # trainable
        code = torch.empty((num_codes,), dtype=torch.uint16)
        self.register_buffer('code', code) # non-trainable
        if bias:
            self.bias = nn.Parameter(torch.empty(out_features))
        else:
            self.register_parameter('bias', None)
    def forward(self, x):
        wetors = self.codebook[self.code]
        if self.deficiency > 0:
        weight = vectors.view(self.in_features, -1)[:, :-self.deficiency]
        else:
            out = torch.matmul(x, weight) + self.bias
        else:
            out = torch.matmul(x, weight)
        return out
```

1184 1185 1186

1213

Listing 1: PyTorch code for the linear layer of ClusComp. All data type is 16-bit.