Reallm: A general framework for LLM compression and fine-tuning

Louis Leconte*

Lisite, Isep, Sorbonne University Math. and Algo. Sciences Lab, Huawei Tech louis.leconte@ens-paris-saclay.fr

Van Minh Nguyen

Math. and Algo. Sciences Lab, Huawei Tech.

Lisa Bedin*

CMAP, Ecole Polytechnique, France lisa.bedin@polytechnique.edu

Eric Moulines

CMAP, Ecole Polytechnique, France

Abstract

We introduce ReALLM, a novel approach for compression and memory-efficient adaptation of pre-trained language models that encompasses most of the post-training quantization and fine-tuning methods for a budget of <4 bits. Pre-trained matrices are decomposed into a high-precision low-rank component and a vector-quantized latent representation (using an autoencoder). During the fine-tuning step, only the low-rank components are updated. Our results show that pre-trained matrices exhibit different patterns. ReALLM adapts the shape of the encoder (small/large embedding, high/low bit VQ, etc.) to each matrix. ReALLM proposes to represent each matrix with a small embedding on b bits and a neural decoder model \mathcal{D}_{ϕ} with its weights on b_{ϕ} bits. The decompression of a matrix requires only one embedding and a single forward pass with the decoder. Our weight-only quantization algorithm yields the best results on language generation tasks (C4 and WikiText-2) for a budget of 3 bits without any training. With a budget of 2 bits, ReALLM achieves state-of-the art performance after fine-tuning on a small calibration dataset.

1 Introduction

Large Language Models (LLMs) based on transformer architectures (Vaswani et al., 2017) have attracted increasing interest, especially with the availability of high-quality, open-source LLMs such as LLaMA (Touvron et al., 2023), Falcon (Almazrouei et al., 2023) and Gemma (Team et al., 2024). These open models offer the advantage that they can be used by end users for inference or local fine-tuning, provided their hardware has sufficient memory for the size of the models. However, "full fine-tuning" — a process that involves updating all previously trained parameters — is still prohibitively expensive for large models. For example, the standard 16-bits fine-tuning of the LLaMA-65B parameter model requires over 780 GB of GPU memory (Dettmers et al., 2023a). This high requirement is due to the need to store both the weights of the model and the states of the optimizer in GPU memory, a need that increases as the size of the LLMs increases.

A common method to mitigate memory constraints is to quantize the model weights, activations, and gradients — to a lower bit precision. Quantization-Aware Training (QAT) is often used in computer vision; see Courbariaux et al. (2015); Liu et al. (2020); Gholami et al. (2022). However, training large language models (LLMs) from scratch is impractical due to high computational cost. Post-training quantization (PTQ) is an efficient compromise (Dettmers et al., 2022; Frantar et al., 2022), which has

^{*}equal contribution

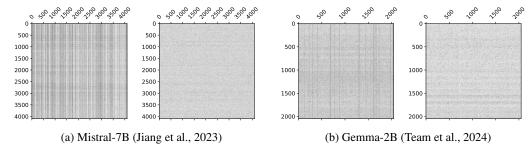


Figure 1: Pre-trained matrix from the first block (left; with "structures"), and pre-trained matrix from the last block (right) for two different models. Stronger vertical patterns appear in the first blocks.

recently attracted much attention (Kim et al., 2023b; Dettmers et al., 2023b; Kim et al., 2023a; Shao et al., 2023). Although most research focuses on scalar quantization (SQ), a few studies investigate LLM compression using vector quantization (VQ) (Tseng et al., 2024; Egiazarian et al., 2024).

In Dettmers et al. (2023a), quantization is effectively combined with the Parameter Efficient Fine-Tuning (PEFT) method, LoRA (Hu et al., 2021), to improve efficiency and practicality in memory-constrained environments. Post-Training Quantization (PTQ) has the potential to be further improved to achieve sub-3 bit quantization (Li et al., 2023; Guo et al., 2023). However, it was found that the weights of the LLM often contain outliers — weights with significantly higher values than others (Kim et al., 2023b; Dettmers et al., 2023b). These outliers pose a considerable challenge for model compression with PTQ and lead to significant quantization errors.

In this paper we present ReALLM - for **Re**sidual Autoencoder **LLM** - a general approach for LLM PTQ and fine-tuning. Pre-trained LLM matrices are decomposed into a 16-bit remainder (low rank, sparse outliers, etc.) and a compressed part, which is fed into a VQ autoencoder (Van Den Oord et al., 2017). In our experiments, we implement a low-rank and quantized decomposition of pre-trained LLM matrices. In this approach, only the low-rank components are fine-tuned (block-wise and end-to-end) while the quantized elements remain static. Our quantization strategy (i.e. the shape of the autoencoder) adapts to the matrix patterns: Our results suggest that some pre-trained LLM matrices exhibit "spatial" patterns (see Figure 1; left) that bear similarities to those in images/videos and allow for highly effective compression (see Figure 3).

Contributions:

- We present ReALLM, a method that uses a novel autoencoder and a residual pipeline to efficiently compress pre-trained LLM matrices;
- We show that state-of-the-art PTQ approaches (Lin et al., 2023; Shao et al., 2023; Tseng et al., 2024; Egiazarian et al., 2024) and fine-tuning methods (Hu et al., 2021; Dettmers et al., 2023a; Guo et al., 2023; Li et al., 2023; Liao and Monz, 2024) are all special cases of ReALLM;
- We propose a preprocessing step that includes scaling and column permutations of matrices to mitigate the quantization errors associated with outliers; We also propose to adapt the general autoencoder scheme to the type of pre-trained matrix patterns.
- Our approach demonstrates that fine-tuning end-to-end with block-wise error reduction leads to the best results reported in the literature for 3 and 2-bit Post-Training Quantization (PTQ).

2 Related works

LLMs adapters. After the introduction of high-performance open-source LLMs and due to the impracticality of "full fine-tuning", several methods of parameter-efficient fine-tuning (PEFT) have emerged, including prefix tuning (Li and Liang, 2021), selective fine-tuning (Guo et al., 2021) and Low Rank Adapter (LoRA). LoRA, introduced in Hu et al. (2021), is a simple but effective fine-tuning method that retains the pre-trained matrices but adds a low-rank component. For a typical pre-trained

matrix W of size 4096×4096 , LoRA introduces two additional matrices of size $4096 \times r$ and $r \times 4096$, where $r \ll 4096$, and tunes only their $2 \times r \times 4096$ parameters. In our work, we use DoRA (Liu et al., 2024) to further improve the fine-tuning by decomposing a weight into its magnitude and direction: $W_{\text{finetune}} = m \frac{W + L_1(L_2)^t}{\|W + L_1(L_2)^t\|_c}$, where W is the frozen pre-trained weight, m is the trainable size vector, (L_1, L_2) are the low-rank (trainable) adapters, and $\|\cdot\|_c$ denotes the Euclidean norm of a matrix over each column. DoRA with the trainable size vector requires little computational effort, but can lead to significant performance improvements (Liu et al., 2024).

Quantization. Current methods for compressing LLMs predominantly use quantization techniques. Early strategies, such as ZeroQuant (Yao et al., 2022) and nuQmm (Park et al., 2022), relied primarily on direct rounding of weights to the nearest quantization level. Later developments improved this approach by handling outliers through quantization to higher bitwidths (Xiao et al., 2023; Dettmers et al., 2022; Kim et al., 2023b; Dettmers et al., 2023b). Methods similar to ReALLM include those that combine quantization with a low-rank decomposition; see e.g. Dettmers et al. (2023a); Guo et al. (2023); Li et al. (2023); Liao and Monz (2024). QLoRA (Dettmers et al., 2023a) combined Parameter Efficient Fine-Tuning (PEFT) and quantization, but added zero-initialised low-rank adapters after quantization. In contrast, Loftq (Li et al., 2023) and LQ-LoRA (Guo et al., 2023) propose to minimize quantization errors by initializing LoRA components with an SVD of the pre-trained weights. As part of this integration, ApiQ (Liao and Monz, 2024) uses gradient descent to optimize both the LoRA components and the quantization parameters for the entire model rather than for each individual layer. Quantization of pre-trained weights facilitates efficient inference on devices with limited memory. To achieve significant computational and energy efficiency, recent studies have combined quantization of weights with activation quantization (Liu et al., 2023; Nrusimha et al., 2024).

Block/Layer-Wise Tuning. GPTQ (Frantar et al., 2022) introduced a higher accuracy strategy using an approximate large-scale solver to minimize the layer-wise quadratic error, which is crucial for low bit-width quantization, as highlighted in Tseng et al. (2024); Egiazarian et al. (2024). Quip# (Tseng et al., 2024) applies random rotations to the pre-trained matrices, segments the resulting matrix into vectors of dimension d=8 and uses optimal lattice quantizers (Viazovska, 2017) to quantize each vector. Due to the random rotation, the distribution of the coefficient vector resembles an isotropic Gaussian distribution, but breaks the inherent dependence between the individual coefficients (see Figure 1). In contrast, AQLM (Egiazarian et al., 2024) uses additive quantization with adaptive codebooks per layer and performs blockwise fine-tuning. Each codebook is first filled with Kmeans (Arthur et al., 2007), and the codewords are optimized to minimize the mean square error caused by the VQ at the output of each block. Quip# and AQLM have achieved stable results (i.e. a single-digit increase in perplexity) in the compression range of 2 bits per parameter.

3 Method

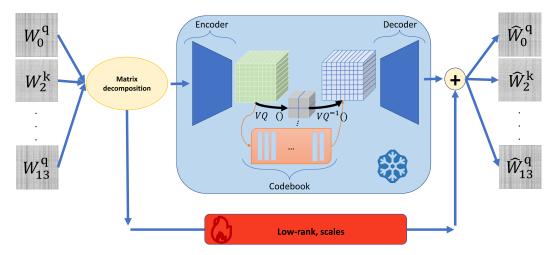


Figure 2: ReALLM; during the fine-tuning step only low-rank and scales are updated

Low-rank/sparse decomposition. Starting from a pre-trained LLM matrix $W \in \mathbb{R}^{p \times q}$, W is decomposed in a first step into a residual component $R \in \mathbb{R}^{p \times q}$ and a quantized matrix Q (which is represented on average with b bits per coordinate). Only the residual matrix is retained with high bit accuracy and further optimized in the fine-tuning phase using a small calibration dataset. Any efficient matrix decomposition can fit into the residual part: butterfly (Dao et al., 2019), sparse outliers (Dettmers et al., 2023b; Lin et al., 2023), etc. In Section 4 we use a low-rank component $R = L_1(L_2)^t$. This structure is analogous to the *data-free* method described in Guo et al. (2023). The aim is to identify Q, L_1 and L_2 that (approximately) solve the following problem:

$$\min_{Q, L_1, L_2} \|W - (Q + L_1(L_2)^t)\|.$$

QLoRA Dettmers et al. (2023a) provides a suboptimal solution for the previously described optimization problem by setting $L_1=0$ and solving $\min_Q\|W-Q\|$. There is no guarantee that the initialization of the low-rank part to zero is optimal. It has been reported that QLoRA, Apiq and Loftq perform better than QLoRA in several language generation benchmarks (Guo et al., 2023; Liao and Monz, 2024; Li et al., 2023).

Mixed-autoencoder configuration. An autoencoder is the composition of an encoding function $\mathcal E$ and a decoding function $\mathcal D$. In ReALLM, $\mathcal E_\psi$ and $\mathcal D_\phi$ are parameterized by neural networks ψ,ϕ and especially $\mathcal E_\psi:\mathbb R^{p\times q}\to\mathbb R^{e_0\times e_1\times e_2}, \mathcal D_\phi:\mathbb R^{e_0\times e_1\times e_2}\to\mathbb R^{p\times q},$ with $e_0e_1e_2\ll pq$. As far as we know, most previous works on quantization of LLMs have focused on applying the same quantization strategy *directly* to the (rotated) pre-trained matrix: i.e. take the embedding dimensions $e_0=p,e_1=q,e_2=1$. Quip# (Tseng et al., 2024) is a special case of ReALLM (with no residual R) where the encoder is assumed to be a (random) rotation matrix $\mathcal E_\psi=U$ and the decoder is assumed to be the inverse $\mathcal D_\phi=U^{-1}$. LQ-LoRA (Guo et al., 2023), Loftq (Li et al., 2023), and ApiQ (Liao and Monz, 2024) are special cases of ReALLM where the encoder and the decoder are defined as the identity matrix.

The approach may not be optimal as some matrices are more challenging to quantize than others (Guo et al., 2023). Specifically, Figure 1 shows that pre-trained LLM matrices can display very different "spatial" patterns. ReALLM adapts the autoencoder to the type and shape of the matrix. When quantizing pre-trained matrices with strong coefficient dependencies, ReALLM is akin to image and video compression techniques that use the implicit neural representation (Chen et al., 2023; Kwan et al., 2024). ReALLM extracts latent representations $\mathcal{E}_{\psi}(W)$ of a set of trained LLM matrices. In the next step, a decoder model is trained to generate the original LLM matrices $\mathcal{D}_{\phi}(\mathcal{E}_{\psi}(W))$ (refer to Figure 2). During the inference phase of an LLM, only the latent embedding $\mathcal{E}_{\psi}(W)$ and the decoder \mathcal{D}_{ϕ} are needed to reconstruct the original weight W, with the exception of the additional low-rank and scale components. We use HNeRV (Chen et al., 2023) to train the autoencoder efficiently. HNeRV (over-)fits a model to the input matrices (i.e. here the pre-trained LLM matrices) with an encoder \mathcal{E}_{ψ} consisting of standard 2D convolutions, and a decoder combining 2D-convNeXt (Liu et al., 2022) and PixelShuffle (Shi et al., 2016).

The decoding process is fast, as HNeRV requires only one network forward operation for decoding. Reallm compression is a combination of a small (w.r.t. input signals) neural decoder model \mathcal{D}_{ϕ} and model compression ($b_{\phi} \ll 16$). HNeRV implements weight pruning (Han et al., 2015), weight quantization (PTQ) and entropy encoding. We go one step further by using a QAT approach: we train the decoder network \mathcal{D}_{ϕ} with convolution kernels quantized to $b_{\phi}=6$ bits during training with the straight-through estimator (Bengio, 2013). For a typical matrix of size 4096×4096 , we train a decoder network with $c=7.2\cdot 10^6$ parameters on $b_{\phi}=6$ bits and an embedding of size $16\times 16\times 16$.

The total bit budget for the given matrix is therefore $\frac{6\cdot(7.2\cdot10^6)+16\cdot(16\cdot16\cdot16\cdot\frac{4096^2}{512^2})}{(4096)^2}=2.82$ bits per coordinate.

Vector Quantization (VQ). An efficient way to store the embedding $\mathcal{E}_{\psi}(W)$ with few bits is VQ. AQLM (Egiazarian et al., 2024) is a special case of Reallm where the latent representation is the matrix W itself. AQLM optimizes multiple codebooks with gradient descent thanks to a calibration dataset. In contrast, for the forward pass, we opted for a *data-free* vector quantization (VQ) method based on Kmeans (Arthur et al., 2007). A given embedding of size $e_0 \times e_1 \times e_2$ is divided into buckets of dimension d. First, we compute scales with NF-normalization (Dettmers et al., 2023a; Guo et al., 2023). The scales are further quantized following the idea of LQ-LoRA, resulting in an additional memory cost of 0.1 bit (Guo et al., 2023). Then we optimize $2^{b \cdot d}$ codewords using Kmeans

clustering on the set of vectors in dimension d to create a codebook. Each vector of dimension d is quantized by the index of the closest element in the codebook (see Figure 2). Consequently, the total number of bits required is $(bd)\frac{e_0e_1e_2}{d}$, i.e. b bits per coordinate. Additional memory is required to store the codebook (namely $16 \cdot d \times 2^{b \cdot d}$ bits). It should be noted that no separate gradient is defined for the quantization operator with the closest element (Van Den Oord et al., 2017). Therefore, during the backward pass, we approximate the gradient similarly to the straight-through estimator (Bengio, 2013) and simply copy the gradients from the decoder input to the encoder output.

Quantization pre-processing. Before using a tensor quantization method, it is important to perform an appropriate scaling. Several parameters (number of blocks, quantile bins, etc.) are chosen to correspond to a given compression ratio. But the presence of outliers (Kim et al., 2023b; Dettmers et al., 2023b) forces the scaling and quantization methods to have a poor compression ratio (Lin et al., 2023; Tseng et al., 2024; Ashkboos et al., 2024). Incoherence processing uses random rotations as a pre-processing step. Although the main purpose of incoherence processing is to reduce the effects of outliers (Tseng et al., 2024; Ashkboos et al., 2024), this technique has a detrimental effect on the structure of the pre-trained matrices within the initial blocks of the LLM (see Figures 1 and 3). This is a serious bottleneck as quantization errors in these initial blocks can propagate throughout the model. As shown in Figure 1, some matrices have no specific patterns and resemble random Gaussian noise interspersed with randomly positioned outliers. To deal with outliers in the latent representation, we suggest rearranging the columns to create some spatial regularity. This strategy aims to find the most effective permutations that cluster outliers. Trukhanov and Soloveychik (2024) has recently elaborated a row/column permutation strategy that summarizes vectors (i.e. sets of rows or columns) with similar norms. In contrast, for ReALLM we propose to permute columns such that neighboring columns are "similar" and not just on the same hypersphere. We develop a basic, yet efficient method for this: first we select a block of size $128 \times q$ in the input tensor of size $e_0 \times e_1 \times e_2$. We start from the first vector, and we search for its closest neighbor in the set of (q-1) vectors (we compute (q-1)scalar products and select the vector that minimizes it). Then, we permute the neighbor vector with the vector in the second position of the block. The process is then iterated; more details are given in Algorithm 1 and Appendix A.3. Note that the memory storage of the permutation is negligible: for a LLM matrix with q = 4096 columns, the permutation requires $q \log(q) = 12 \times 4096$ additional bits for each block of size 128×4096 , hence the memory overhead is about 0.09 bits per coordinate.

Algorithm 1: permutation function

```
 \begin{array}{c|c} \textbf{Input :} \textbf{Matrix } w \text{ of size } 128 \times q \text{ ;} \\ \textbf{1 for } j = 0, \ldots, q-1 \text{ do} \\ \textbf{2} & | column_j = w[:,j] \text{ ;} \\ \textbf{3} & | indx_j = get\_index\_nn(column_j, w[:,j+1:q]) /* \text{ get the nearest neighbor index of current } column_j \text{ , among the rest of un-permuted columns} \\ & | w[:,j+1:q] & | */ \\ \textbf{4} & | \text{Permute } w[:,j+1] \text{ and } w[:,indx_j]; \\ \textbf{5} & | \text{Save the inverse of the permutation index in } inv\_permut; \\ \textbf{6} & \text{end} \\ & | \textbf{Output :} w, inv\_permut \\ \end{array}
```

Reallm: a new LLM format. LLM standard formats represent LLM weights as a set of matrices encoded on 16 bits. Scalar quantization approaches (Frantar et al., 2022; Dettmers et al., 2023a) represent any matrix of size $p \times q$ with $b \cdot pq$ bits for a budget of b bits. Vector quantization (VQ) methods (Egiazarian et al., 2024; Tseng et al., 2024) represent any matrix of size $p \times q$ with a smaller matrix of size $p \times \frac{q}{d}$ with $b \cdot d$ bits for a budget of b bits and a vector dimension d. Reallem goes one step further and proposes to represent each matrix of size $p \times q$ with a small embedding of size $e_0 \times e_1 \times e_2$ on b bits and a neural decoder model \mathcal{D}_ϕ with c parameters on b_ϕ bits. Figure 2 illustrates the most important innovation of Reallem: LLMs are no longer represented by a set of matrices, but as a combination of embeddings and a single neural decoder model. Reallem learns a single model for a specific family of basic models (e.g. LlaMAs, Gemma). If a specific weight matrix is needed for a specific LLM, one must take its embedding and perform only one single forward pass with the decoder \mathcal{D}_ϕ . This speeds up the decoding step compared to diffusion-based approaches (Wang et al.,

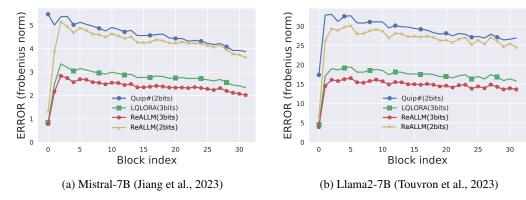


Figure 3: Reconstruction (Frobenius norm) error for layer of type "Q" for all blocks. Quip# (Tseng et al., 2024) does not take advantage of the structures in the first blocks.

2024; Soro et al., 2024). Note that for ReALLM a decoder model has to be trained on LLM matrices, but this learning step is done once and for all. Additionally, the more we train and overfit, the better ReALLM becomes.

The set of hyper-parameters for ReallM are: r the rank, (e_0,e_1,e_2) the shape of the latent representation, (b,d) the number of bits and the bucket dimension in the VQ, and (c,b_ϕ) the number of parameters and the number of bits of the decoder. We have conducted extensive experiments to find suitable configurations; however, we were unable to test configurations with a large decoder size. For e.g., for small embeddings $(e_0e_1e_2<1024)$ and a total budget of 3 bits for a single LLaMA2-7B model, the decoder model in ReallM has $c=3.5\cdot 10^9$ parameters trained on $b_\phi=6$ bits. Our GPU is unable to accommodate multiple LLM matrices in memory for ReallM training, typically with size $n\times n$; n>4096. Therefore, we test ReallM on a set of 512×512 patches extracted from pre-trained LLM matrices, and we use the HNeRV Chen et al. (2023) autoencoder model. For more details on the practical aspect of decoder training, see Appendix A.2.

We have experimentally discovered two sets of optimal combinations of hyperparameters that depend on the type and shape of the pre-trained matrix. Some pre-trained matrices, especially those closer to the input tokens, compress better with small latent representations $(e_0e_1e_2<1024)$ in high bit precision (b>8) and (relatively) large decoders $(c>4\cdot 10^6)$. Other pre-trained matrices (usually deeper in the LLM) compress better with very large embeddings $(e_0>\frac{p}{4},e_1>\frac{q}{4},e_2\in[1,2])$ with low bit precision $(b\ll8)$ and (relatively) small decoders $(c\ll10^6)$. In Figure 3 ReALLM achieves the lowest Frobenius norm quantization error. We perform ablation experiments with this metric to decouple the effects of VQ and permutation preprocessing of ReALLM on the final performance. For example, in block 8 (Mistral-7b; left panel), the error for scalar quantization (SQ; used in Dettmers et al. (2023a); Guo et al. (2023)) is 2.96. This error decreases with VQ to 2.68 and with permutation further to 2.54, while permutation alone (i.e. with SQ) leads to an error of 2.88. Quip# rotates the matrices randomly, causing all patterns in the initial blocks to be lost.

4 Experimental validation

We test ReALLM on the LLaMA-2 (Touvron et al., 2023) family models (with 7 and 13 billions parameters). We compare our method with other quantization approaches for a budget of 3 and 2 bits per coordinate. We partially reused code from the implementations of LQ-LoRA², AQLM ³ and HNeRV⁴. On an Nvidia A40 GPU (with 46GB memory), the entire computation (PTQ + fine tuning) takes about 90 hours for a LLaMA2-7B model.

Language Generation Tasks. For continual language modeling, we train on a single partition of the C4 (Raffel et al., 2020) dataset for half an epoch and use a sequence length of 4096 for training

²https://github.com/HanGuo97/lq-lora/tree/main

³https://github.com/Vahe1994/AQLM

⁴https://github.com/haochen-rye/HNeRV

Algorithm 2: Pseudo-code for ReALLM with block-wise and end-to-end fine-tuning

Input: Number of end-to-end fine-tuning steps T, Number of block-wise fine-tuning steps K, Number of blocks n, Shape of the latent space (e_0, e_1, e_2) , Number of weights in the decoder c, Number of bits for the decoder weights b_ϕ , Number of VQ bits per dimension b, VQ dimension d, Rank r;

```
1 Initialize
         Get pre-trained matrices \{W^q, W^k, W^v, W^o, W^{gate}, W^{up}, W^{down}\} for all n blocks;
 2
 3 end
    /* Block-wise fine-tuning
                                                                                                                                   */
 4 for j = 0, ..., n-1 do
         B_j = \{W^q, W^k, W^v, W^o, W^{gate}, W^{up}, W^{down}\}[block = j];
         output_j = forward\_pass(B_j) / * get \textit{non-quantized} output
 6
         for l \in \{q, k, v, o, gate, up, down\} do
 7
              \begin{split} L1_j^l, L2_j^l &= svd\_decomposition(W_j^l, rank = r) \;; \\ W_j^l &= W_j^l - L1_j^l(L2_j^l)^t \;; \\ \mathcal{E}_{\psi}(W_j^l), \mathcal{D}_{\phi_j^l} &= autoencoder(W_j^l, e_0, e_1, e_2, c, b_{\phi}) \; /* \; \text{latent representation} \end{split}
 8
10
              \mathcal{E}_{\psi}(W_j^l), inv\_permut_j^l = permute(\mathcal{E}_{\psi}(W_j^l)) \; / * \; \text{with Algorithm 1}
11
              \mathcal{E}_{\psi}(W_j^l) = normalize(\mathcal{E}_{\psi}(W_j^l)) /* with NF-normalization (Dettmers
12
                   et al., 2023a; Guo et al., 2023)
              codebook_i^l = Kmeans(\mathcal{E}_{\psi}(W_i^l), b, d);
13
              codes_{j}^{l} = get\_index\_nn(\mathcal{E}_{\psi}(\mathring{W}_{j}^{l}), codebook_{j}^{l}) /* get nearest neighbor index
14
             \begin{split} W_j^l \leftarrow \{codes_j^l, codebook_j^l, \mathcal{D}_{\phi_j^l}, inv\_permut_j^l, L1_j^l, L2_j^l\}; \\ dora_j^l = DoRA(W_j^l, L1_j^l, L2_j^l) \ /* \ \text{get DoRA scale} \end{split}
15
16
17
         end
         dora\_quantized\_output_j = forward\_pass\_quantized(\{dora_j^l, L1_j^l, L2_j^l, W_j^l\}_{l \ge 0})
18
              /* get output after quantization and DoRA
         L_i = ||output_i - dora\_quantized\_output_i||^2;
19
         for k = 0, ..., K - 1 do
20
         Optimize \{dora_j^l, L1_j^l, L2_j^l\}_{l\geq 0} with gradient descent to minimize L_j;
21
22
         end
23 end
    /* End-to-end fine-tuning
                                                                                                                                   */
24 for t = 0, \dots, T - 1 do
         Optimize \{dora_i^l, L1_i^l, L2_i^l\}_{l,j>0} with gradient descent;
25
26 end
```

only. Note that the native context length for LLaMA-2 (Touvron et al., 2023) is 4096, while it is 2048 for LLaMA-1. Consequently, in the literature LLaMA-2 models are evaluated with token sequences of size 2048 (all except (Egiazarian et al., 2024) follow this rule). Therefore, we use a sequence length of size 2048 for both WikiText-2 (Merity et al., 2016) and C4 evaluation.

Our main baselines are LQ-LoRA (Guo et al., 2023), Quip# (Tseng et al., 2024), and AQLM (Egiazarian et al., 2024). However, we also report the performance of popular quantization approaches GPTQ (Frantar et al., 2022), AWQ (Lin et al., 2023), Omniquant (Shao et al., 2023), as well as the performance of recent work ApiQ (Liao and Monz, 2024) and QuaRot (Ashkboos et al., 2024). In the results below, we present the target bits per parameter that takes into account quantized weights and include parameters kept in high precision (head layer, scales, codebooks, permutations in 16 bits, and low-rank matrices in 8 bits precision) similarly to the related work. The exact bit budget is detailed in Table 5 in the Appendix.

In our experiments, following Dettmers et al. (2023a); Guo et al. (2023), we take a DoRA (Liu et al., 2024) rank of r = 64 (unless otherwise specified), we set the decoder bit precision to $b_{\phi} = 6$, and

we adjust the size of the latent representation (e_0, e_1, e_2) depending on the block index (tested from (4096, 4096, 1) to (16, 16, 16)), and we have tested several VQ in dimension d=2 or d=4. The VQ-autoencoder is trained with cosine scheduler with a maximum learning rate of 0.001 for 2000 epochs. Then we (optionally) tune the low-rank components block-wise with a batch of size 32 and a step size of $1 \cdot e^{-5}$. The end-to-end fine-tuning is run with batches of size 1, and a learning rate of $2 \cdot e^{-5}$. As far as we know, we have also developed the first VQ code (available in the supplementary material) that makes efficient use of PyTorch's "torch dispatch" functionality (Ansel et al., 2024), which is known to be as fast as dedicated CUDA kernels (Guo et al., 2023). This allows us to overload PyTorch operations to perform just-in-time dequantization.

In Tables 1 and 2 we evaluate the perplexity of Reallem on the respective validation datasets of C4 and WikiText-2 for a single run. During fine-tuning (on a single partition of the C4 dataset), we only update the DoRA components (scales and low-rank matrices). For each dataset, we provide three sets of results in Table 1: Perplexity without any fine-tuning (only low-rank and VQ autoencoder decomposition), perplexity with only block-wise fine-tuning, and perplexities with end-to-end fine-tuning (in addition to the block-wise fine-tuning process). Our *data-free* version of Reallem (no

Table 1: Perplexity (↓) on the validation dataset for LLaMA2-7B, with a sequence length of 2048

Method	#bits	$\operatorname{rank} r$	bucket d	C4 (↓)	WikiText-2 (↓)
ReALLM (no fine-tuning)	3	32	2	7.78	6.21
ReALLM (block-wise)	3	32	2	7.56	6.01
ReALLM (40% training)	3	32	2	7.31	5.80
ReALLM (full training)	3	32	2	7.29	5.79
ReALLM (no fine-tuning)	3	64	2	7.72	6.10
ReALLM (block-wise)	3	64	2	7.51	5.92
ReALLM (40% training)	3	64	2	7.30	5.78
ReALLM (full training)	3	64	2	7.27	5.77
ReALLM (no fine-tuning)	2	64	2	45.96	51.74
ReALLM (block-wise 50 epochs)	2	64	2	18.61	16.95
ReALLM (block-wise 200 epochs)	2	64	2	10.11	8.31
ReALLM (40% training)	2	64	2	8.56	6.95
ReALLM (full training)	2	64	2	8.47	6.91
ReALLM (no fine-tuning)	2	64	4	41.02	40.85
ReALLM (block-wise 50 epochs)	2	64	4	15.74	12.08
ReALLM (40% training)	2	64	4	8.36	6.74
ReALLM (full training)	2	64	4	8.28	6.69

fine-tuning; see Table 1) achieves state-of-the-art metrics for 3 bit quantization. However, for a budget of 2 bits, quantization errors are larger, and our results show that fine-tuning (both block-wise and end-to-end) is needed to further improve performance. This result is in line with the PTQ literature (Frantar et al., 2022; Egiazarian et al., 2024). Table 1 also shows that reducing the rank from r = 64to r=32 has minimal effect on the final perplexity result, while halving the number of parameters that need to be tuned. Moreover, a larger VQ dimension d=4 instead of d=2 leads to better results. Note that increasing d comes at an additional storage cost (as explained in Section 3, $16 \cdot d \times 2^{b \cdot d}$ bits are needed to store the codebook). Additional results for other models are available in the Appendix. In Table 2 we compare ReALLM with end-to-end fine-tuning, and the best performing PTQ approaches. All the methods cited in Table 2 also uses a calibration dataset. It is interesting to note that ReALLM with 2 bits bridges the gap with the famous GPTQ (Frantar et al., 2022) method on 3 bits for the LLaMA2-13B. One major difference between ReALLM and Quip# (Tseng et al., 2024) is that the quantized weights are kept frozen during all the fine-tuning process in ReALLM. As a consequence, we can store a single version of the quantized weight, and fine-tune several versions of the learnable parameters (i.e. DoRA scales and low-rank matrices) for several fine-tuning tasks. On the contrary Quip# updates all the weights (in 16 bits precision) during the layer-wise fine-tuning. This does not only slow down the PTQ process (as gradients must be store for all weights in the given block), but it also means Quip# has to store learnable vectors and also quantized weights for each fine-tuning task.

Table 2: Perplexity (\downarrow) on the validation dataset for LLaMA2-7B and LLaMA2-13B, with a sequence length of 2048

Method	Number of bits	C4 (↓)		WikiText-2 (↓)	
		7B	13B	7B	13B
LLaMA2 (Touvron et al., 2023)	16	6.97	6.46	5.47	4.48
GPTQ (Frantar et al., 2022)	3	7.89	7.00	6.29	5.42
AWQ (Lin et al., 2023)	3	7.84	6.94	6.24	5.32
Omniquant (Shao et al., 2023)	3	7.75	6.98	6.03	5.28
LQ-LoRA (Guo et al., 2023)	3	7.88	_	6.48	_
LoftQ (Li et al., 2023)	3	_	_	5.63	5.13
ApiQ[PTQ] (Liao and Monz, 2024)	3	7.84	6.88	6.19	5.18
Quip# (Tseng et al., 2024)	3	7.32	6.72	5.79	5.10
QuaRot[A16W3] (Ashkboos et al., 2024)	3	_	_	6.09	5.37
ReALLM	3	7.27	6.69	5.77	5.14
LoftQ (Li et al., 2023)	2	l –	_	7.85	7.69
ApiQ (Liao and Monz, 2024)	2	_	_	7.46	6.29
Quip# (Tseng et al., 2024)	2	8.35	7.45	6.66	5.74
AQLM (Egiazarian et al., 2024)	2	8.56	7.51	6.64	5.65
ReALLM	2	8.28	7.50	6.69	5.72

Table 3: Accuracy (↑) in LM Eval (acc, not acc_norm).

Method	Size	#bits	ARC-challenge	ARC-easy	PiQA	Winogrande	Average
LLaMA-2	7B	16	43.52	76.26	78.07	69.22	66.77
AQLM (Egiazarian et al., 2024)	7B	2	33.55	62.79	73.54	64.61	58.62
Quip# (Tseng et al., 2024)	7B	2	34.63	64.60	75.12	64.89	59.81
ReALLM	7B	2	35.15	68.56	75.73	66.46	61.47
LLaMA-2	13B	16	48.32	78.48	80.01	72.13	69.74
AQLM (Egiazarian et al., 2024)	13B	3	43.63	73.51	77.78	67.56	65.62
Quip# (Tseng et al., 2024)	13B	3	44.02	72.45	78.40	69.13	66.00
ReALLM	13B	3	47.01	75.96	78.67	70.96	68.15

Zero-Shot Tasks. Following HuggingFace's Open LLM Leaderboard⁵, and the literature (Frantar et al., 2022; Guo et al., 2023), we also measure zero-shot accuracy on ARC (Clark et al., 2018), PiQA (Tata and Patel, 2003), and Winogrande (Sakaguchi et al., 2021), via the LM Evalaluation Harness (Gao et al., 2021). We report results in Table 3, and compute the average on the 4 mentioned tasks. For all LLM sizes, ReALLM provides a notable advantage (between 0.5 and 3 points of accuracy improvement) with respect to AQLM (Egiazarian et al., 2024) and Quip# (Tseng et al., 2024). Interestingly, the LLaMA2-13B model compressed on 3 bits with ReALLM performs better than the standard LLaMA-2-7B model (16 bits) on the zero-shot tasks.

5 Conclusion

We present ReALLM, a weight-only PTQ method that achieves state-of-the-art results on LLMs at 2, and 3 bits budget. Our (low-rank) fine-tuning approach enables one to fine-tune language models with 13 billions parameters on a *single* GPU with less than 40 GB of RAM.

Large context sequence lengths result in large KV-cache memory consumption during inference, and PTQ is a promising approach for compressing KV-cache activations (Hooper et al., 2024; Ashkboos et al., 2024). Concurrently to our work, Trukhanov and Soloveychik (2024) propose a quantization method based on permutations of rows from K and V matrices. We are currently studying how to adapt ReALLM to KV-cache quantization, and how to combine it with activation quantization.

⁵https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard

6 Societal impact

This paper presents work whose goal is to advance the field of LLM compression and fine-tuning. There are many potential societal consequences of our work, in particular malicious usage of LLMs for spams or language generation on edge devices. However, this negative societal impact is not limited to ReALLM, but to the field of LLM in general.

References

- Almazrouei, E., Alobeidli, H., Alshamsi, A., Cappelli, A., Cojocaru, R., Debbah, M., Goffinet, É., Hesslow, D., Launay, J., Malartic, Q., et al. (2023). The falcon series of open language models. *arXiv* preprint arXiv:2311.16867.
- Ansel, J., Yang, E., He, H., Gimelshein, N., Jain, A., Voznesensky, M., Bao, B., Bell, P., Berard, D., Burovski, E., Chauhan, G., Chourdia, A., Constable, W., Desmaison, A., DeVito, Z., Ellison, E., Feng, W., Gong, J., Gschwind, M., Hirsh, B., Huang, S., Kalambarkar, K., Kirsch, L., Lazos, M., Lezcano, M., Liang, Y., Liang, J., Lu, Y., Luk, C., Maher, B., Pan, Y., Puhrsch, C., Reso, M., Saroufim, M., Siraichi, M. Y., Suk, H., Suo, M., Tillet, P., Wang, E., Wang, X., Wen, W., Zhang, S., Zhao, X., Zhou, K., Zou, R., Mathews, A., Chanan, G., Wu, P., and Chintala, S. (2024). PyTorch 2: Faster Machine Learning Through Dynamic Python Bytecode Transformation and Graph Compilation. In 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2 (ASPLOS '24). ACM.
- Arthur, D., Vassilvitskii, S., et al. (2007). k-means++: The advantages of careful seeding. In *Soda*, volume 7, pages 1027–1035.
- Ashkboos, S., Mohtashami, A., Croci, M. L., Li, B., Jaggi, M., Alistarh, D., Hoefler, T., and Hensman, J. (2024). Quarot: Outlier-free 4-bit inference in rotated llms. *arXiv preprint arXiv:2404.00456*.
- Bengio, Y. (2013). Estimating or propagating gradients through stochastic neurons. *arXiv* preprint *arXiv*:1305.2982.
- Chen, H., Gwilliam, M., Lim, S.-N., and Shrivastava, A. (2023). Hnerv: A hybrid neural representation for videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10270–10279.
- Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. (2018). Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv* preprint *arXiv*:1803.05457.
- Courbariaux, M., Bengio, Y., and David, J.-P. (2015). Binaryconnect: Training deep neural networks with binary weights during propagations. *Advances in neural information processing systems*, 28.
- Dao, T., Gu, A., Eichhorn, M., Rudra, A., and Ré, C. (2019). Learning fast algorithms for linear transforms using butterfly factorizations. In *International conference on machine learning*, pages 1517–1527. PMLR.
- Dettmers, T., Lewis, M., Belkada, Y., and Zettlemoyer, L. (2022). Gpt3. int8 (): 8-bit matrix multiplication for transformers at scale. *Advances in Neural Information Processing Systems*, 35:30318–30332.
- Dettmers, T., Pagnoni, A., Holtzman, A., and Zettlemoyer, L. (2023a). Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36.
- Dettmers, T., Svirschevski, R. A., Egiazarian, V., Kuznedelev, D., Frantar, E., Ashkboos, S., Borzunov, A., Hoefler, T., and Alistarh, D. (2023b). Spqr: A sparse-quantized representation for near-lossless llm weight compression. In *The Twelfth International Conference on Learning Representations*.
- Egiazarian, V., Panferov, A., Kuznedelev, D., Frantar, E., Babenko, A., and Alistarh, D. (2024). Extreme compression of large language models via additive quantization. *arXiv* preprint *arXiv*:2401.06118.
- Frantar, E., Ashkboos, S., Hoefler, T., and Alistarh, D. (2022). Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv* preprint arXiv:2210.17323.
- Gao, L., Tow, J., Abbasi, B., Biderman, S., Black, S., DiPofi, A., Foster, C., Golding, L., Hsu, J., Le Noac'h, A., Li, H., McDonell, K., Muennighoff, N., Ociepa, C., Phang, J., Reynolds, L., Schoelkopf, H., Skowron, A., Sutawika, L., Tang, E., Thite, A., Wang, B., Wang, K., and Zou, A. (2021). A framework for few-shot language model evaluation.

- Gholami, A., Kim, S., Dong, Z., Yao, Z., Mahoney, M. W., and Keutzer, K. (2022). A survey of quantization methods for efficient neural network inference. In *Low-Power Computer Vision*, pages 291–326. Chapman and Hall/CRC.
- Guo, D., Rush, A. M., and Kim, Y. (2021). Parameter-efficient transfer learning with diff pruning. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 4884–4896.
- Guo, H., Greengard, P., Xing, E., and Kim, Y. (2023). Lq-lora: Low-rank plus quantized matrix decomposition for efficient language model finetuning. In *The Twelfth International Conference on Learning Representations*.
- Han, S., Pool, J., Tran, J., and Dally, W. (2015). Learning both weights and connections for efficient neural network. *Advances in neural information processing systems*, 28.
- Hooper, C., Kim, S., Mohammadzadeh, H., Mahoney, M. W., Shao, Y. S., Keutzer, K., and Gholami, A. (2024). Kvquant: Towards 10 million context length llm inference with kv cache quantization. *arXiv* preprint arXiv:2401.18079.
- Hu, E. J., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W., et al. (2021). Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., Casas, D. d. l., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., et al. (2023). Mistral 7b. arXiv preprint arXiv:2310.06825.
- Kim, J., Lee, J. H., Kim, S., Park, J., Yoo, K. M., Kwon, S. J., and Lee, D. (2023a). Memory-efficient fine-tuning of compressed large language models via sub-4-bit integer quantization. arXiv preprint arXiv:2305.14152.
- Kim, S., Hooper, C., Gholami, A., Dong, Z., Li, X., Shen, S., Mahoney, M. W., and Keutzer, K. (2023b). Squeezellm: Dense-and-sparse quantization. *arXiv preprint arXiv:2306.07629*.
- Kwan, H. M., Gao, G., Zhang, F., Gower, A., and Bull, D. (2024). Hinery: Video compression with hierarchical encoding-based neural representation. *Advances in Neural Information Processing Systems*, 36.
- Li, X. L. and Liang, P. (2021). Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597.
- Li, Y., Yu, Y., Liang, C., Karampatziakis, N., He, P., Chen, W., and Zhao, T. (2023). Loftq: Lora-fine-tuning-aware quantization for large language models. In *The Twelfth International Conference on Learning Representations*.
- Liao, B. and Monz, C. (2024). Apiq: Finetuning of 2-bit quantized large language model. *arXiv* preprint arXiv:2402.05147.
- Lin, J., Tang, J., Tang, H., Yang, S., Dang, X., and Han, S. (2023). Awq: Activation-aware weight quantization for llm compression and acceleration. *arXiv* preprint arXiv:2306.00978.
- Liu, J., Gong, R., Wei, X., Dong, Z., Cai, J., and Zhuang, B. (2023). Qllm: Accurate and efficient low-bitwidth quantization for large language models. In *The Twelfth International Conference on Learning Representations*.
- Liu, S.-Y., Wang, C.-Y., Yin, H., Molchanov, P., Wang, Y.-C. F., Cheng, K.-T., and Chen, M.-H. (2024). Dora: Weight-decomposed low-rank adaptation. *arXiv preprint arXiv:2402.09353*.
- Liu, Z., Mao, H., Wu, C.-Y., Feichtenhofer, C., Darrell, T., and Xie, S. (2022). A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11976–11986.

- Liu, Z., Shen, Z., Savvides, M., and Cheng, K.-T. (2020). Reactnet: Towards precise binary neural network with generalized activation functions. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16*, pages 143–159. Springer.
- Merity, S., Xiong, C., Bradbury, J., and Socher, R. (2016). Pointer sentinel mixture models. In *International Conference on Learning Representations*.
- Nrusimha, A., Mishra, M., Wang, N., Alistarh, D., Panda, R., and Kim, Y. (2024). Mitigating the impact of outlier channels for language model quantization with activation regularization. arXiv preprint arXiv:2404.03605.
- Park, G., Park, B., Kim, M., Lee, S., Kim, J., Kwon, B., Kwon, S. J., Kim, B., Lee, Y., and Lee, D. (2022). Lut-gemm: Quantized matrix multiplication based on luts for efficient inference in large-scale generative language models. *arXiv preprint arXiv:2206.09557*.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.
- Sakaguchi, K., Bras, R. L., Bhagavatula, C., and Choi, Y. (2021). Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106.
- Shao, W., Chen, M., Zhang, Z., Xu, P., Zhao, L., Li, Z., Zhang, K., Gao, P., Qiao, Y., and Luo, P. (2023). Omniquant: Omnidirectionally calibrated quantization for large language models. In *The Twelfth International Conference on Learning Representations*.
- Shi, W., Caballero, J., Huszár, F., Totz, J., Aitken, A. P., Bishop, R., Rueckert, D., and Wang, Z. (2016). Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1874–1883.
- Soro, B., Andreis, B., Lee, H., Chong, S., Hutter, F., and Hwang, S. J. (2024). Diffusion-based neural network weights generation. *arXiv* preprint arXiv:2402.18153.
- Tata, S. and Patel, J. M. (2003). Piqa: An algebra for querying protein data sets. In 15th International Conference on Scientific and Statistical Database Management, 2003., pages 141–150. IEEE.
- Team, G., Mesnard, T., Hardin, C., Dadashi, R., Bhupatiraju, S., Pathak, S., Sifre, L., Rivière, M., Kale, M. S., Love, J., et al. (2024). Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al. (2023). Llama: Open and efficient foundation language models. *arXiv* preprint arXiv:2302.13971.
- Trukhanov, N. and Soloveychik, I. (2024). Accurate block quantization in llms with outliers. *arXiv* preprint arXiv:2403.20137.
- Tseng, A., Chee, J., Sun, Q., Kuleshov, V., and De Sa, C. (2024). Quip#: Even better Ilm quantization with hadamard incoherence and lattice codebooks. *arXiv preprint arXiv:2402.04396*.
- Van Den Oord, A., Vinyals, O., et al. (2017). Neural discrete representation learning. *Advances in neural information processing systems*, 30.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. Advances in neural information processing systems, 30.
- Viazovska, M. S. (2017). The sphere packing problem in dimension 8. Annals of mathematics, pages 991–1015.
- Wang, K., Xu, Z., Zhou, Y., Zang, Z., Darrell, T., Liu, Z., and You, Y. (2024). Neural network diffusion. *arXiv preprint arXiv:2402.13144*.

- Xiao, G., Lin, J., Seznec, M., Wu, H., Demouth, J., and Han, S. (2023). Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*, pages 38087–38099. PMLR.
- Yao, Z., Yazdani Aminabadi, R., Zhang, M., Wu, X., Li, C., and He, Y. (2022). Zeroquant: Efficient and affordable post-training quantization for large-scale transformers. *Advances in Neural Information Processing Systems*, 35:27168–27183.

A Appendix / supplemental material

A.1 Structures in pre-trained matrices

Interestingly, the blocks that show some visual structures in LLaMA and Mistral models are not the same for Gemma LLMs. For instance in Figure 4, we can see that Gemma2b (Team et al., 2024)'s matrices keep some internal patterns in all blocks, not only at the very first blocks. Note this has no negative impact on ReALLM, as the shape of the encoder is experimentally adapted to each block.

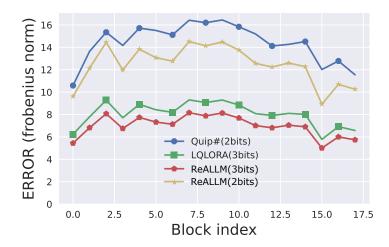


Figure 4: Reconstruction (Frobenius norm) error for layer of type "Q" for all blocks of Gemma2b LLM.

A.2 Autoencoder computational limitations

Our GPU can not directly work on LLM pre-trained matrices with large sizes (typically of shape 4096×4096). Instead, we choose to split each pre-trained matrix into a set of 64 "patches" of shapes 512×512 , and we learn the decoder on the set of matches rather than on the big initial matrix. During the inference time, when de-quantizing a LLM matrix, we reconstruct each patch (in parallel) and we concatenate the patches together. This step of concatenation has a minimal impact on the final time complexity of our method. In Table 4, we present ablation experiment results on the type of decoder weight (only) quantization. We performed a quantization aware training approach, i.e. directly optimizing weight quantized on b_{ϕ} bits using straight through estimator Bengio (2013). We also tested a post training quantization method where the weight of the decoder are quantized with a round to nearest (RTN) approache, at the end of the decoder training steps. We

Table 4: Reconstruction (Frobenius norm) error for layer of type "Q" inside the first block of Mistral-7b model, for patches of size 512×512 using a constant embedding size of $(e_0, e_1, e_2) = (16, 16, 16)$, and a varying quantization strategy (during the decoder training, i.e. QAT, or after the training, i.e. PTQ).

Error	# parameters c ($\times 10^6$)	b_{ϕ}	bit budget	quantization
0.84	_	_	3	NF3(Guo et al., 2023)
1.78	7.2	6	2.82	PTQ
1.19	5.4	7	2.44	PTQ
1.61	7.7	5	2.32	QAT
1.24	4.5	8	2.21	QAT
0.69	7.2	6	2.82	QAT

vary the number of parameters c, and the bit precision b_{ϕ} of the decoder to target a total bit cost below 3 bits per coordinate. This experiment show two different results: first, the influence on the

quantization performance of the number of decoder parameters c and their respective bit precision b_{ϕ} is not straightforward. Second, under the same parameters (number of parameters and bits), QAT gives better performance than the respective PTQ approach. Furthermore, for a reduced number of bits (2.82 vs 3), ReALLM yields a smaller quantization error compared to the scalar quantization NF3 (Dettmers et al., 2023a; Guo et al., 2023).

Table 5: Comparison of several LLM format for m matrices of size $p \times q$, and a budget of b bits per coordinate. ReALLM uses a decoder model with c parameters trained on b_{ϕ} bits, and a rank r.

Method	LoRA	VQ only (like AQLM)	ReALLM
Matrix representation	$(p \times q) \cdot 16$	$(p \times \frac{q}{d}) \cdot b \cdot d$ $2^{bd} \cdot d \cdot 16$	$(e_0 \times \frac{e_1}{d} \times e_2) \cdot b \cdot d$ $2^{bd} \cdot d \cdot 16$
Codebook	_	$2^{bd} \cdot d \cdot 16$	$2^{b\overline{d}}\cdot d\cdot 16$
Decoder	_	_	cb_{ϕ}
Low-rank	$(2 \times r \times \min(p,q)) \cdot 16$	_	$(2 \times r \times \min(p,q)) \cdot 16$
Total bit cost	$16(pq + 2r\min(p,q)) \cdot m$	$(bpq + 2^{bd+4}d) \cdot m$	$cb_{\phi} + 32r\min(p,q) + m(16d2^{bd} + e_0e_1e_2b)$

Table 6: Quantization and fine-tuning approaches as particular case of ReALLM (with a rank r, and a budget of b bits for VQ in dimension d) for a matrix of size $p \times q$.

_					
Method	\mid rank r	Autoencoder	Latent (e_0, e_1, e_2)	VQ dim. (d)	VQ bits (b)
LoRA (Hu et al., 2021)	64	None	(p, q, 1)	1	16
GPTQ (Frantar et al., 2022)	0	None	(p,q,1)	1	4
QLoRA (Dettmers et al., 2023a)	64	None	(p,q,1)	1	4
LQ-LoRA (Guo et al., 2023)	64	None	(p,q,1)	1	3
Quip# (Tseng et al., 2024)	0	Rotation matrix	(p,q,1)	8	2
AQLM (Egiazarian et al., 2024)	0	None	(p,q,1)	8	2
ReALLM	64	Trainable	(e_0, e_1, e_2)	4	2

A.3 Permutations

In ReALLM, we compute permutations on sets of vectors in dimension 128. We could work with smaller blocks, but it induces more memory dedicated to the permutation storage (one permutation for each block).

We start from the first vector (i.e. the first column of the initial matrix shrunk to a dimension d=128), and we search for its closest neighbor in the set of (q-1) vectors (we compute (q-1) scalar products and select the vector that minimizes it). Then, we permute the neighbor vector with the vector in the second position of the block. We then focus on the second vector, and search for its closest neighbor in the set of (q-2) vectors. The process is then iterated. Details are given in Algorithm 1.

A.4 Broader impacts and Safeguards

Our computing unit seriously restricts the size of the decoder models we can train. We are not able to train one decoder model for a given LLM, but we work layer-wise and train a single decoder model for all patches extracted from the given layer. This layer-wise training forms the main limitation of Reallm w.r.t. standard post-training quantization methods, such as round to nearest (RTN).

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

Table 7: Perplexity (\downarrow) on the validation dataset for LLaMA2-13B, with a sequence length of 2048

Method	#bits	$\operatorname{rank} r$	bucket d	C4 (↓)	WikiText-2 (↓)
ReALLM (no fine-tuning)	3	64	2	6.91	5.27
ReALLM (30% training)	3	64	2	6.69	5.14
ReALLM (no fine-tuning)	2	64	4	10.36	8.15
ReALLM (10% training)	2	64	4	7.59	5.99