

---

# AAAC: Activation-Aware Adaptive Codebooks for 4-bit LLM Weight Quantization

---

Beshr IslamBouli<sup>\*1</sup> David Jin<sup>\*1</sup>

## Abstract

Post-training weight-only quantization to 4 bits is widely used to reduce the memory and compute costs of large language model inference. Existing PTQ methods, such as AWQ and GPTQ, improve how weights are mapped onto a fixed 4-bit grid through scaling, clipping, or error compensation. To further improve accuracy, methods such as OmniQuant and QuIP# use gradient-assisted algorithms at the cost of hours of quantization time. In this work, we propose AAAC (Activation-Aware Adaptive Codebooks), a lightweight method for 4-bit LLM weight quantization. AAAC replaces the fixed scalar codebook used in standard quantization with two small learned scalar codebooks (64 bytes) per layer. Each group of weights selects the codebook that minimizes activation-weighted reconstruction error, encoding the choice in the unused sign bit of the group’s positive scale and adding zero storage overhead. AAAC completes in 3–30 minutes on a single GPU, and adds no memory beyond the model itself. We evaluate against AWQ, GPTQ, IF4, GPTVQ, OmniQuant, SqueezeLLM, and QuIP# across model families. AAAC outperforms baselines at orders-of-magnitude less quantization time.

## 1. Introduction

Large language models (LLMs) are widely deployed under memory and compute constraints that make weight quantization essential (Touvron et al., 2023; Qwen et al., 2025). Weight-only 4-bit quantization has emerged as a practical inference setting: it reduces model storage while retaining acceptable quality for many models (Dettmers et al., 2023).

---

<sup>\*</sup>Equal contribution <sup>1</sup>Massachusetts Institute of Technology, Cambridge, MA, USA. Correspondence to: Beshr IslamBouli <beshr@mit.edu>, David Jin <jindavid@mit.edu>.

Presented at the 43<sup>rd</sup> International Conference on Machine Learning Workshop on Resource-Adaptive Foundation Model Inference (AdaptFM), Seoul, South Korea, 2026. Copyright 2026 by the author(s).

Most post-training quantization (PTQ) methods improve how weights are mapped onto a fixed 4-bit grid. In group-wise quantization, each group of weights shares a scale  $s$ , and each normalized weight  $\tilde{w} = w/s$  is rounded to a discrete code under a format-defined table  $\mathcal{T}$ . The same table is then used during dequantization:

$$c = Q(w; s), \quad \hat{w} = Q^{-1}(c; s) = \mathcal{T}[c] \cdot s. \quad (1)$$

For INT4,  $\mathcal{T}$  is an integer grid; for NVFP4,  $\mathcal{T}$  is the FP4 grid. Once the format is chosen, this scalar grid is usually treated as fixed. Methods such as AWQ and GPTQ therefore improve quantization around this fixed grid, using channel scaling (Lin et al., 2024), clipping, or error compensation (Frantar et al., 2023), while leaving the scalar reconstruction grid itself unchanged. We argue that this fixed-grid assumption leaves quality on the table. The dequantized values, not the discrete codes alone, determine the weights actually used in W4A16 and W4A8 inference. A 4-bit format fixes the number of stored codes and the group-wise scale structure, but it need not fix the best scalar grid for a given layer. This suggests a complementary PTQ direction: instead of only asking how to map weights onto a fixed FP4 or INT4 table, learn adaptive scalar codebooks that better match the weight distribution and the activation statistics of the layer.

We propose **Activation-Aware Adaptive Codebooks (AAAC)**, a lightweight method for 4-bit LLM weight quantization. AAAC replaces the single fixed table  $\mathcal{T}$  with two learned scalar codebooks per layer,  $\mathcal{T}_0$  and  $\mathcal{T}_1$ . Each group of weights selects the codebook that minimizes activation-weighted reconstruction error, and weights in that group are stored as ordinary 4-bit indices into the selected codebook. AWQ and GPTQ improve how weights are mapped onto a fixed 4-bit table; AAAC improves the table used for that mapping.

AAAC distinguishes itself from prior PTQ methods on three axes. **Negligible overhead:** AAAC adds only two BF16 codebooks (64 bytes) per layer. The per-group selection bit is stored in the unused sign of the always-positive scale, adding zero storage. **No gradients:** codebooks are learned via activation-weighted  $k$ -means on forward-pass activations alone—no backpropagation, no Hessian or Fisher computation. **No memory blowup:** AAAC processes each layer

independently using only its activations and codebook statistics, so peak memory stays close to the model size itself. Across Llama and Qwen models in both NVFP4 and INT4 settings, we evaluate against AWQ, GPTQ, IF4, GPTVQ, OmniQuant, SqueezeLLM, and QuIP#. AAAC outperforms baselines at orders-of-magnitude less quantization time.

## 2. Related Work

Post-training quantization methods for LLMs can be divided into gradient-free approaches, which use only forward-pass calibration data, and gradient-assisted approaches, which incur larger quantization times and/or peak memory overhead.

**Gradient-free PTQ.** GPTQ (Frantar et al., 2023) quantizes weights column-by-column using inverse-Hessian error compensation. AWQ (Lin et al., 2024) protects salient weights through activation-aware channel scaling. IF4 (Cook et al., 2026) selects per-group between FP4 and INT4 codebooks based on MSE. All three complete in minutes with peak memory close to the original model size.

**Gradient-assisted PTQ.** OmniQuant (Shao et al., 2024) learns per-channel clipping thresholds and equivalent transforms via block-wise backpropagation through the rounding operation, requiring 1–16 hours for 7B–70B models.

SqueezeLLM (Kim et al., 2024) learns per-channel non-uniform lookup tables via Fisher-weighted k-means, where the Fisher diagonal is computed by backpropagating through the full model. While the algorithm runs in 11–80 minutes, the Fisher step drives peak memory to 2–4× model size (33–292 GB for 7B–65B).

QuIP# (Tseng et al., 2024) applies a Randomized Hadamard Transform for incoherence processing, then rounds to a fixed E8 lattice codebook via BlockLDLQ. Without fine-tuning this is gradient-free and completes in under 10 GPU-hours for 70B. An optional fine-tuning stage adds backpropagation through the full model, raising the cost to roughly 100 GPU-hours.

GPTVQ (van Baalen et al., 2025) extends GPTQ to vector quantization, learning per-group codebooks via Hessian-weighted EM followed by gradient descent on the layer-wise reconstruction objective to refine codebook entries. Quantization takes 0.5–11 hours for 7B–70B models.

Our experiments show that gradient-assisted methods consistently outperform gradient-free methods, suggesting that gradients—with their associated quantization time and memory overhead—are necessary to close the gap to full-precision quality. AAAC challenges this conclusion: it is a lightweight, gradient-free method that learns only two codebooks per layer (64 bytes overhead) via activation-weighted k-means. Despite this simplicity, AAAC outperforms all

gradient-free methods and matches or exceeds gradient-assisted methods, with 3–30 minutes of quantization time and peak memory no larger than the original model size.

## 3. Method

### 3.1. Background: Group-wise Quantization

Consider a linear layer computing  $\mathbf{y} = \mathbf{x}\mathbf{W}^\top$  with  $\mathbf{W} \in \mathbb{R}^{N \times K}$ . In group-wise quantization, each row of  $\mathbf{W}$  is partitioned into groups of  $g$  contiguous elements. Each group shares a scale  $s > 0$ , and each weight is represented by a discrete code.

**Quantization.** A scale  $s$  is computed from the group’s weight statistics. Each weight  $w$  is normalized by  $s$ , rounded to the nearest entry in the format-defined table  $\mathcal{T}$ , and stored as a  $b$ -bit code:

$$\tilde{w} = w/s, \quad c = \arg \min_i |\tilde{w} - \mathcal{T}[i]|. \quad (2)$$

**Dequantization.** During inference, the stored code is looked up in the same table and multiplied by the scale:

$$\hat{w} = \mathcal{T}[c] \cdot s. \quad (3)$$

Thus the table  $\mathcal{T}$  plays two roles: it defines the nearest-neighbor code assignment during quantization, and it defines the reconstruction value used during dequantization.

It is useful to distinguish the stored code from the reconstructed value. For any scalar table  $\mathcal{A}$ , define the nearest-entry reconstruction operator

$$\text{recon}_{\mathcal{A}}(\tilde{w}) = \mathcal{A} \left[ \arg \min_i |\tilde{w} - \mathcal{A}[i]| \right]. \quad (4)$$

This operator first finds the nearest entry of  $\mathcal{A}$  to the normalized weight  $\tilde{w}$ , then returns the corresponding table value. Under the base format, the reconstructed normalized weight is  $\text{recon}_{\mathcal{T}}(\tilde{w})$ , and the dequantized weight after quantization is  $\text{recon}_{\mathcal{T}}(\tilde{w}) \cdot s$ .

The table  $\mathcal{T}$  is fixed by the format. For NVFP4,  $\mathcal{T}$  contains the signed E2M1 values used by the FP4 format, such as  $\{-6, -4, -3, \dots, 0, \dots, 3, 4, 6\}$ . For INT4,  $\mathcal{T}$  denotes the signed or zero-point-adjusted integer grid used by the implementation, expressed in normalized units. Thus both NVFP4 and INT4 can be written as group-wise scalar quantizers with a finite code alphabet and a per-group scale.

**Key observation.** Existing PTQ methods such as AWQ and GPTQ primarily improve quantization around the fixed table  $\mathcal{T}$ . They change the weight distribution, scales, or compensation procedure so that the fixed table produces better reconstructed values. AAAC targets a different degree of freedom: it learns the scalar table itself.

### 3.2. AAAC: Activation-Aware Adaptive Codebooks

AAAC replaces the single fixed table  $\mathcal{T}$  with two learned scalar codebooks  $\mathcal{T}_0$  and  $\mathcal{T}_1$ . Each learned table has the same number of entries as the base 4-bit format, but its entries are optimized from calibration data rather than fixed by the format. A group of  $S$  consecutive weights selects one of the two tables. The selected table defines the nearest-entry reconstruction used for that group.

**Granularity parameters.** AAAC has two granularity parameters. The *scale group size*  $g$  is fixed by the quantization format and determines how many weights share one scale:  $g=16$  for NVFP4 and  $g=128$  for INT4. The *selection group size*  $S$  determines how many weights share one codebook choice. Thus, every  $S$  weights store one selection bit  $\sigma$ , and every  $g$  weights store one scale.

When  $S=g$ , the selection group coincides with the scale group: each scale needs exactly one codebook-selection bit. Since the scale is always positive, its sign bit is unused, and we store  $\sigma$  there with no additional selection storage. When  $S < g$ , the selection metadata costs  $1/S$  bits per weight. Smaller  $S$  gives finer codebook adaptation. We use  $S=g=16$  for NVFP4 and evaluate  $S \in \{16, 128\}$  for INT4.

**Modified quantization and dequantization.** For a weight  $w$  in selection group  $j$ , let  $\tilde{w} = w/s$ . AAAC selects table  $\mathcal{T}_{\sigma_j}$  and reconstructs the normalized weight using the nearest-entry reconstruction operator:

$$\hat{w} = \text{recon}_{\mathcal{T}_{\sigma_j}}(\tilde{w}) \cdot s. \quad (5)$$

Equivalently, the stored 4-bit code is the nearest-entry index under the selected table:

$$c = \arg \min_i |\tilde{w} - \mathcal{T}_{\sigma_j}[i]|. \quad (6)$$

Thus AAAC preserves the scalar 4-bit code width, group-wise scale structure, and packed representation, but it does not require the code assignment to match round-to-nearest under the original table  $\mathcal{T}$ . Its improvement comes from learning better scalar reconstruction grids and choosing, per group, which grid to use.

**Activation importance.** For a linear layer  $\mathbf{y} = \mathbf{x}\mathbf{W}^\top$ , a perturbation  $\delta\mathbf{W}$  produces output error  $\delta\mathbf{y} = \mathbf{x}\delta\mathbf{W}^\top$ . The expected squared output error contributed by column  $k$  of  $\delta\mathbf{W}$  is proportional to  $\mathbb{E}[x_k^2] \|\delta\mathbf{W}_{:,k}\|^2$ . Summing over a calibration set  $\mathbf{X} \in \mathbb{R}^{T \times K}$  gives the per-column importance

$$I_k = \sum_{t=1}^T X_{t,k}^2. \quad (7)$$

Minimizing  $\sum_k I_k \|\delta\mathbf{W}_{:,k}\|^2$  is therefore a diagonal proxy for minimizing calibration output error. Columns with larger

activation energy receive larger weight in the reconstruction objective.

**Table selection.** Consider a selection group  $j$  containing normalized weights  $\{\tilde{w}_{j,1}, \dots, \tilde{w}_{j,S}\}$  and activation importances  $\{I_{j,1}, \dots, I_{j,S}\}$ . AAAC assigns the group to the table that minimizes activation-weighted reconstruction error:

$$\sigma_j = \arg \min_{r \in \{0,1\}} \sum_{k=1}^S I_{j,k} (\tilde{w}_{j,k} - \text{recon}_{\mathcal{T}_r}(\tilde{w}_{j,k}))^2. \quad (8)$$

The table choice and code assignment are coupled: each candidate table defines both the reconstruction values and the nearest-entry regions used to encode the weights in the group.

**Storage.** Each learned table contains  $|\mathcal{T}|$  values stored in BF16, for a total of  $2|\mathcal{T}| \cdot 2$  bytes per layer, which is at most 64 bytes. For a 7B model with  $\sim 224$  linear layers, this is about 14 KB, negligible compared with the model weights. The selection bits  $\{\sigma_j\}$  require no additional storage when  $S=g$  if they are stored in the unused sign bit of each positive scale, and require  $1/S$  bits per weight when  $S < g$ .

### 3.3. Learning Algorithm

We now describe how the tables  $\mathcal{T}_0$  and  $\mathcal{T}_1$  are obtained from calibration data. The two tables are learned per layer using a small calibration set. Learning alternates between assigning selection groups to tables and refining the table entries. Algorithm 1 provides pseudocode.

**Initialization.** Let  $M = |\mathcal{T}|$  be the number of entries in each table, and let  $\text{Quantile}_q(\{\tilde{w}\})$  denote the  $q$ -quantile of the normalized weight distribution. We initialize  $\mathcal{T}_0$  using  $M$  evenly spaced quantiles spanning the full range:

$$\mathcal{T}_0[i] = \text{Quantile}_{i/(M-1)}(\{\tilde{w}\}), \quad i = 0, \dots, M-1. \quad (9)$$

We initialize  $\mathcal{T}_1$  using  $M$  evenly spaced quantiles starting at a half-step offset and ending at the maximum:

$$\mathcal{T}_1[i] = \text{Quantile}_{\delta+(1-\delta)i/(M-1)}(\{\tilde{w}\}), \quad \delta = \frac{1}{2(M-1)}. \quad (10)$$

This shifts the second grid forward by half a quantile spacing, giving the two tables distinct but nearby starting points. Results are insensitive to the exact offset.

**Alternating optimization.** Learning alternates between two steps for  $n_{\text{outer}}$  iterations. In the *assignment step*, each selection group is assigned to the table that minimizes its activation-weighted reconstruction error, as in Eq. 8. In the *table update step*, each table is refined by  $n_{\text{inner}}$  iterations of activation-weighted scalar  $k$ -means. For each table,

**Algorithm 1** AAAC: Learning activation-aware adaptive codebooks for one layer

**Require:** Weight matrix  $\mathbf{W} \in \mathbb{R}^{N \times K}$ , calibration activations  $\mathbf{X} \in \mathbb{R}^{T \times K}$

**Require:** Group sizes  $g$  (scale),  $S$  (selection); iterations  $n_{\text{outer}}, n_{\text{inner}}$

- 1: Compute scales  $s$  per group of  $g$  weights
- 2: Normalize weights:  $\tilde{\mathbf{W}} \leftarrow \mathbf{W}/s$
- 3: Compute activation importance:  $I_k \leftarrow \sum_t X_{t,k}^2$  for each column  $k$
- 4: Initialize  $\mathcal{T}_0$  from full-range quantiles of  $\tilde{\mathbf{W}}$  {Eq. 9}
- 5: Initialize  $\mathcal{T}_1$  from half-shifted quantiles of  $\tilde{\mathbf{W}}$  {Eq. 10}
- 6: **for**  $t = 1$  to  $n_{\text{outer}}$  **do**
- 7:   **// Assignment step**
- 8:   **for** each selection group  $j$  **do**
- 9:      $\sigma_j \leftarrow \arg \min_{r \in \{0,1\}}$
- 10:      $\sum_{k=1}^S I_{j,k} (\tilde{w}_{j,k} - \text{recon}_{\mathcal{T}_r}(\tilde{w}_{j,k}))^2$
- 11:   **end for**
- 12:   **// Table update step**
- 13:   **for**  $r \in \{0, 1\}$  **do**
- 14:     Let  $\mathcal{W}_r$  be the multiset of pairs  $(\tilde{w}_{j,k}, I_{j,k})$  with  $\sigma_j = r$
- 15:     **for**  $\ell = 1$  to  $n_{\text{inner}}$  **do**
- 16:       **for** each entry  $i$  in  $\mathcal{T}_r$  **do**
- 17:         Let  $\mathcal{V}_{r,i} \subseteq \mathcal{W}_r$  be the Voronoi cell of entry  $i$
- 18:          $\mathcal{T}_r[i] \leftarrow \frac{\sum_{(\tilde{w}, I) \in \mathcal{V}_{r,i}} I \tilde{w}}{\sum_{(\tilde{w}, I) \in \mathcal{V}_{r,i}} I}$
- 19:         If the denominator is zero, keep the previous value of  $\mathcal{T}_r[i]$
- 20:       **end for**
- 21:     **end for**
- 22:     Sort  $\mathcal{T}_r$  in ascending order
- 23:   **end for**
- 24: **end for**
- 25: Round  $\mathcal{T}_0, \mathcal{T}_1$  to BF16 precision
- 26: Recompute  $\{\sigma_j\}$  using Eq. 8
- 27: Pack 4-bit codes as nearest-entry indices into  $\mathcal{T}_{\sigma_j}$
- 28: Store  $\{\sigma_j\}$  in sign bits of  $s$  if  $S=g$ , or separately if  $S < g$   $\mathcal{T}_0, \mathcal{T}_1$ , packed 4-bit codes, scales, and  $\{\sigma_j\}$

normalized weights assigned to that table are assigned to their nearest table entry, and each entry is updated to the activation-weighted centroid of its Voronoi cell:

$$\mathcal{T}_r[i] \leftarrow \frac{\sum_{(\tilde{w}, I) \in \mathcal{V}_{r,i}} I \tilde{w}}{\sum_{(\tilde{w}, I) \in \mathcal{V}_{r,i}} I}. \quad (11)$$

Here,  $\mathcal{V}_{r,i}$  is the set of normalized weights assigned to entry  $i$  of table  $\mathcal{T}_r$ . If a Voronoi cell is empty, we keep the entry at its previous value. After the final inner update, the entries of each table are sorted in ascending order. The outer loop matters more than the inner loop: reassigning groups to tables has a larger effect on quality than refining table

entries within a fixed assignment. We find empirically that a single outer iteration already achieves 82% of the total improvement, so the overall cost of the procedure is dominated by the assignment step rather than the inner  $k$ -means refinement.

**Computational cost.** Each outer iteration requires two forward-only passes over the layer’s weights—one for group assignment, one for  $k$ -means updates—with no gradients, no Hessian-vector products, and no backpropagation. The algorithm processes layers independently, so peak GPU memory is dominated by the model itself plus a per-layer activation buffer. This enables in-memory calibration of a 27B model in roughly 10 minutes ( $n_{\text{outer}}=3, n_{\text{inner}}=10$ ).

## 4. Experiments

### 4.1. Setup

**Models.** We evaluate on models spanning Llama (3.2-1B, 3.2-3B, 3.1-8B) (Grattafiori et al., 2024), Qwen2.5 (0.5B, 7B, 14B, 32B) (Qwen et al., 2025), and Qwen3.5 (2B, 4B, 9B, 27B) (Qwen Team, 2026). To compare with numbers reported by prior methods, we evaluate on LLaMA-1 (7B, 13B) and LLaMA-2 (7B, 13B) (Touvron et al., 2023).

**Evaluation.** We report perplexity on WikiText-2 (Merity et al., 2016) and C4 (Raffel et al., 2023). For WikiText-2, we follow the GPTQ protocol: the full test set is tokenized without filtering and split into non-overlapping 2048-token segments, with per-token cross-entropy computed via manual logit shifting (Frantar et al., 2023). For C4, we use the validation split with up to 256 segments. We additionally evaluate on downstream tasks using lm-evaluation-harness (Gao et al., 2024).

**Gap recovery.** To summarize the overall effectiveness of each method, we report *gap recovery*: the percentage of the quantization-induced degradation that a method eliminates relative to the round-to-nearest baseline. For perplexity (lower is better), gap recovery is defined as  $(\text{PPL}_{\text{RTN}} - \text{PPL}_{\text{method}}) / (\text{PPL}_{\text{RTN}} - \text{PPL}_{\text{full}}) \times 100\%$ , where  $\text{PPL}_{\text{RTN}}$  is the round-to-nearest perplexity and  $\text{PPL}_{\text{full}}$  is the full-precision baseline. For accuracy (higher is better), the analogous formula is  $(\text{Acc}_{\text{method}} - \text{Acc}_{\text{RTN}}) / (\text{Acc}_{\text{full}} - \text{Acc}_{\text{RTN}}) \times 100\%$ . A recovery of 0% means no improvement over RTN; 100% means full-precision quality is restored.

**AAAC configurations.** A smaller selection granularity  $S$  allows AAAC to better adapt to local weight structure within each layer. The NVFP4 memory layout groups weights into blocks of 16 with shared FP8 scales, and this  $g=16$  setting is ideal for AAAC: 16 weights is small enough for fine-grained codebook selection, and the selection bit can be

Table 1.  $g=16$  W4A16 perplexity on WikiText-2 ( $\downarrow$ ). AAAC uses  $S=g=16$  (zero overhead). Best quantized result in **bold**.

Model	BF16	RTN	IF4	AWQ	AAAC
Llama-3.2-1B	9.71	10.68	10.61	10.46	<b>10.21</b>
Llama-3.2-3B	7.77	8.19	8.13	8.11	<b>7.99</b>
Llama-3.1-8B	6.19	6.56	6.51	6.51	<b>6.39</b>
Qwen3.5-2B	12.06	13.21	12.70	12.56	<b>12.37</b>
Qwen3.5-4B	9.42	10.25	10.21	10.16	<b>9.92</b>
Qwen3.5-9B	8.51	9.11	9.38	8.91	<b>8.56</b>
Qwen3.5-27B	6.80	7.02	6.92	7.03	<b>6.88</b>
<i>Average</i>	<i>8.64</i>	<i>9.29</i>	<i>9.21</i>	<i>9.11</i>	<i><b>8.90</b></i>
<i>Recovery</i>	—	—	12.3%	28.1%	<b>59.2%</b>

stored in the scale sign at zero overhead. IF4 uses the same  $g=16$  layout, and AWQ can be straightforwardly adapted to it. We report  $g=16$  results alongside both methods. To compare with prior methods that use the standard INT4 layout with  $g=128$ , we also report  $S=g=128$  (zero overhead). However,  $S=128$  is too coarse for effective codebook selection, so we additionally report  $S=16$ , adding a negligible  $1/16=0.0625$  bpw overhead. We use  $n_{\text{outer}}=3$ ,  $n_{\text{inner}}=10$  with 4 calibration sequences from C4. We used a H100 GPU.

**Baselines and reproducibility.** We compare against round-to-nearest (RTN), IF4 (Cook et al., 2026), AWQ (Lin et al., 2024), GPTQ (Frantar et al., 2023), GPTVQ (van Baalen et al., 2025), OmniQuant (Shao et al., 2024), SqueezeLLM (Kim et al., 2024), and QuIP# (Tseng et al., 2024). GPTQ numbers are taken from Lin et al. (2024). We implement IF4 following Cook et al. (2026): for each group, both FP4 and scaled INT4 quantizations are evaluated, and the option with lower MSE is selected. For AWQ at  $g=128$ , we use the official AutoAWQ package; for  $g=16$ , we implement their per-channel scaling algorithm ourselves as AutoAWQ does not support FP4-based quantization. GPTVQ, OmniQuant, SqueezeLLM, and QuIP# numbers are taken from their respective papers; we verify that our RTN and full-precision baselines match those reported, ensuring consistent evaluation protocols.

## 4.2. Comparison with Gradient-Free Methods

### 4.2.1. GROUP SIZE 16 (NVFP4 MEMORY LAYOUT)

The NVFP4 memory layout stores weights in groups of 16 with shared FP8 scales, and several recent methods build on this layout. IF4 (Cook et al., 2026) selects per-group between the native FP4 and a scaled INT4 codebook. AWQ (Lin et al., 2024), originally designed for INT4, can be straightforwardly adapted to this layout. AAAC uses  $S=g=16$ , storing the selection bit in the scale sign at zero overhead. Table 1 presents WikiText-2 perplexity across 7 Llama and Qwen3.5 models.

Table 2.  $g=128$  W4A16 WikiText-2 perplexity ( $\downarrow$ ) vs. AWQ. AWQ uses the official AutoAWQ package. **Bold**: best result per model.

Model	FP16	RTN	AWQ	AAAC	
				$S=128$	$S=16$
Llama-3.2-3B	7.81	8.49	8.26	8.16	<b>8.10</b>
Llama-3.1-8B	6.24	6.82	6.65	6.68	<b>6.56</b>
Qwen2.5-0.5B	13.07	15.54	15.01	14.47	<b>14.27</b>
Qwen2.5-7B	6.85	7.23	7.09	7.13	<b>7.07</b>
Qwen2.5-14B	5.29	5.78	5.69	5.65	<b>5.60</b>
<i>Average</i>	<i>7.85</i>	<i>8.77</i>	<i>8.54</i>	<i>8.42</i>	<i><b>8.32</b></i>
<i>Recovery</i>	—	—	25.0%	38.0%	<b>48.9%</b>

Table 3.  $g=128$  W4A16 WikiText-2 perplexity ( $\downarrow$ ) vs. GPTQ. GPTQ numbers from Lin et al. (2024). **Bold**: best result per model.

Model	FP16	RTN	GPTQ	AAAC	
				$S=128$	$S=16$
LLaMA-1-7B	5.68	5.96	5.85	5.82	<b>5.76</b>
LLaMA-1-13B	5.09	5.25	5.23	5.21	<b>5.17</b>
LLaMA-2-7B	5.47	5.72	5.69	5.62	<b>5.58</b>
LLaMA-2-13B	4.88	4.98	4.98	4.98	<b>4.94</b>
<i>Average</i>	<i>5.28</i>	<i>5.48</i>	<i>5.44</i>	<i>5.41</i>	<i><b>5.36</b></i>
<i>Recovery</i>	—	—	20.0%	35.0%	<b>60.0%</b>

### 4.2.2. GROUP SIZE 128 (INT4 LAYOUT)

To compare with prior methods that use the standard INT4 layout with  $g=128$ , we evaluate AAAC at  $S=g=128$  (zero overhead) and  $S=16$  (+0.0625 bpw). Table 2 compares AAAC against AWQ on Llama-3 and Qwen2.5 models, using the official AutoAWQ package. With  $S=g=128$  (zero overhead), AAAC already outperforms AWQ on the majority of models. The remaining cases where AWQ leads are attributable to the coarse selection granularity: with only one codebook choice per 128 weights, the method cannot adapt to local variations within each group. With  $S=16$ , adding only 0.0625 bpw, AAAC consistently outperforms AWQ on all models. Table 3 compares AAAC against GPTQ on LLaMA-1 and LLaMA-2 models. GPTQ numbers are taken from Lin et al. (2024). AAAC at  $S=16$  outperforms GPTQ on all four models. Even at  $S=128$  (zero overhead), AAAC matches or outperforms GPTQ on all four models.

## 4.3. Comparison with Gradient-Assisted Methods

As discussed in Section 2, gradient-assisted PTQ methods achieve strong results but at significant cost: OmniQuant requires 1–16 hours of block-wise backpropagation, SqueezeLLM requires  $2\text{--}4\times$  model size in peak memory for Fisher computation, GPTVQ requires 0.5–11 hours with gradient-based codebook refinement, and QuIP#-FT requires up to 100 GPU-hours of end-to-end fine-tuning. Even QuIP# without fine-tuning takes up to 10 GPU-hours.

Table 4.  $g=128$  W4A16 WikiText-2 perplexity ( $\downarrow$ ) vs. gradient-assisted methods. AAAC uses  $S=16$ . Best quantized result per row in bold.

Model	FP16	RTN	OmniQ	SqLLM	GPTVQ	QuIP#	QuIP#-FT	AAAC
LLaMA-1-7B	5.68	5.96	5.77	5.77	5.96	5.83	<b>5.76</b>	<b>5.76</b>
LLaMA-1-13B	5.09	5.25	<b>5.17</b>	<b>5.17</b>	<b>5.15</b>	5.20	<b>5.17</b>	<b>5.17</b>
LLaMA-2-7B	5.47	5.72	5.58	5.57	5.62	5.66	<b>5.56</b>	5.58
LLaMA-2-13B	4.88	4.98	4.95	4.96	4.97	5.00	4.95	<b>4.94</b>
<i>Average</i>	5.28	5.48	5.37	5.37	5.43	5.42	5.36	<b>5.36</b>
<i>Recovery</i>	—	—	55.0%	55.0%	25.0%	30.0%	60.0%	<b>60.0%</b>

Table 5. Combining AWQ and AAAC on WikiText-2 ( $\downarrow$ ). AAAC uses  $S=g$  (zero overhead). Best quantized result per row in bold.

$g$	Model	Full	RTN	AWQ	AAAC	AWQ+AAAC
128	LLaMA-1-7B	5.68	5.96	5.81	5.82	<b>5.80</b>
128	LLaMA-1-13B	5.09	5.25	5.20	5.21	<b>5.18</b>
128	LLaMA-2-7B	5.47	5.72	<b>5.62</b>	<b>5.62</b>	<b>5.62</b>
128	LLaMA-2-13B	4.88	4.98	<b>4.97</b>	4.98	4.98
16	Qwen3.5-4B	9.42	10.25	10.16	9.92	<b>9.52</b>
16	Qwen3.5-9B	8.51	9.11	8.91	<b>8.56</b>	8.67
16	Qwen3.5-27B	6.80	7.02	7.03	6.88	<b>6.80</b>
<i>Average</i>		6.55	6.90	6.81	6.71	<b>6.65</b>
<i>Recovery</i>		—	—	24.2%	53.3%	<b>70.5%</b>

Table 4 compares AAAC against all four methods.

#### 4.4. Combining AAAC with AWQ

AWQ modifies the weights before quantization via channel scaling, while AAAC learns the codebook used during quantization. Since they operate on different components, the two methods can be combined: AWQ first scales the weights, then AAAC learns codebooks on the scaled distribution. Table 5 shows that AWQ+AAAC recovers 70.5% of the quantization gap on average, compared to 53.3% for AAAC alone.

#### 4.5. Downstream Task Evaluation

To complement the perplexity results, we evaluate on 8 downstream tasks (BBH, MuSR, MMLU, GPQA, HellaSwag, ARC-C, PIQA, MATH) using Qwen3.5-9B and Qwen3.5-27B at  $g=16$ . AAAC achieves the highest average accuracy among all quantized methods (64.21%) compared with the BF16 baseline (64.59%). Full per-task results are in Appendix C.

#### 4.6. Additional Analysis

**W4A8 inference.** Although AAAC uses activations to weight importance during calibration, it is not sensitive to activation precision at inference time. The same quantized model can serve both W4A16 and W4A8 inference without re-calibration. We verify this in Appendix D, where the same  $g=128$  models are evaluated with FP8 activations; AAAC’s improvements carry over.

**Sensitivity to hyperparameters.** AAAC has few tunable factors: the selection group size  $S$ , the iteration budget, and the calibration set size. At  $g=16$ ,  $S=g=16$  is the natural choice at zero overhead. At  $g=128$ ,  $S=128$  (zero overhead) already outperforms AWQ on the majority of models (Tables 2–3), while reducing  $S$  to 16 (+0.0625 bpw) consistently improves further. The algorithm is largely insensitive to the iteration budget: a fast configuration ( $n_{\text{outer}}=3$ ,  $n_{\text{inner}}=10$ , under 2 minutes on an H100 for a 7B model) matches or trails something bigger like ( $n_{\text{outer}}=10$ ,  $n_{\text{inner}}=30$ ) by 0.01 perplexity. Calibration data requirements are similarly minimal: in our experience, 1 sequence (2,048 tokens) yields perplexity within 0.01 of the result with 4 or more sequences. We use 4 sequences as a default.

## 5. Conclusion

We have presented AAAC (Activation-Aware Adaptive Codebooks), a lightweight post-training quantization method that learns two small scalar codebooks per layer (64 bytes) via activation-weighted k-means on a small calibration set. By encoding the per-group codebook selection in the unused sign bit of the scale factor, AAAC adds zero storage cost when the selection and scale group sizes match. The method requires no gradients and adds no peak memory beyond the original model size. Despite this simplicity, AAAC consistently outperforms all gradient-free baselines and matches gradient-assisted methods that require hours of optimization and higher memory. We believe that learned codebooks at this minimal scale represent a practical direction for quantized language models.

## Impact Statement

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

## Acknowledgements

We thank Han Guo, Tarushii Goel, and Yoon Kim for helpful discussions.

## References

- Cook, J., Lee, H. S., Le, K., Guo, J., Traverso, G., Chandrakasan, A. P., and Han, S. Adaptive block-scaled data types, 2026. URL <https://arxiv.org/abs/2603.28765>.
- Dettmers, T., Svirschevski, R., Egiazarian, V., Kuznedelev, D., Frantar, E., Ashkboos, S., Borzunov, A., Hoefler, T., and Alistarh, D. Spqr: A sparse-quantized representation for near-lossless llm weight compression, 2023. URL <https://arxiv.org/abs/2306.03078>.
- Frantar, E., Ashkboos, S., Hoefler, T., and Alistarh, D. Gptq: Accurate post-training quantization for generative pre-trained transformers, 2023. URL <https://arxiv.org/abs/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., McDonnell, 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. The language model evaluation harness, 07 2024. URL <https://zenodo.org/records/12608602>.
- Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., Yang, A., Fan, A., Goyal, A., Hartshorn, A., Yang, A., Mitra, A., Sravankumar, A., Korenev, A., Hinsvark, A., Rao, A., Zhang, A., Rodriguez, A., Gregerson, A., Spataru, A., Roziere, B., Biron, B., Tang, B., Chern, B., Caucheteux, C., Nayak, C., Bi, C., Marra, C., McConnell, C., Keller, C., Touret, C., Wu, C., Wong, C., Ferrer, C. C., Nikolaidis, C., Allonsius, D., Song, D., Pintz, D., Livshits, D., Wyatt, D., Esiobu, D., Choudhary, D., Mahajan, D., Garcia-Olano, D., Perino, D., Hupkes, D., Lacomkin, E., AlBadawy, E., Lobanova, E., Dinan, E., Smith, E. M., Radenovic, F., Guzmán, F., Zhang, F., Synnaeve, G., Lee, G., Anderson, G. L., Thattai, G., Nail, G., Mialon, G., Pang, G., Cucurell, G., Nguyen, H., Korevaar, H., Xu, H., Touvron, H., Zarov, I., Ibarra, I. A., Kloumann, I., Misra, I., Evtimov, I., Zhang, J., Copet, J., Lee, J., Geffert, J., Vranes, J., Park, J., Mahadeokar, J., Shah, J., van der Linde, J., Billock, J., Hong, J., Lee, J., Fu, J., Chi, J., Huang, J., Liu, J., Wang, J., Yu, J., Bitton, J., Spisak, J., Park, J., Rocca, J., Johnstun, J., Saxe, J., Jia, J., Alwala, K. V., Prasad, K., Upasani, K., Plawiak, K., Li, K., Heafield, K., Stone, K., El-Arini, K., Iyer, K., Malik, K., Chiu, K., Bhalla, K., Lakhota, K., Rantala-Yeary, L., van der Maaten, L., Chen, L., Tan, L., Jenkins, L., Martin, L., Madaan, L., Malo, L., Blecher, L., Landzaat, L., de Oliveira, L., Muzzi, M., Pasupuleti, M., Singh, M., Paluri, M., Kardas, M., Tsimpoukelli, M., Oldham, M., Rita, M., Pavlova, M., Kambadur, M., Lewis, M., Si, M., Singh, M. K., Hassan, M., Goyal, N., Torabi, N., Bashlykov, N., Bogoychev, N., Chatterji, N., Zhang, N., Duchenne, O., Çelebi, O., Alrassy, P., Zhang, P., Li, P., Vasic, P., Weng, P., Bhargava, P., Dubal, P., Krishnan, P., Koura, P. S., Xu, P., He, Q., Dong, Q., Srinivasan, R., Ganapathy, R., Calderer, R., Cabral, R. S., Stojnic, R., Raileanu, R., Maheswari, R., Girdhar, R., Patel, R., Sauvestre, R., Polidoro, R., Sumbaly, R., Taylor, R., Silva, R., Hou, R., Wang, R., Hosseini, S., Chennabasappa, S., Singh, S., Bell, S., Kim, S. S., Edunov, S., Nie, S., Narang, S., Raparthy, S., Shen, S., Wan, S., Bhosale, S., Zhang, S., Vandenhende, S., Batra, S., Whitman, S., Sootla, S., Collot, S., Gururangan, S., Borodinsky, S., Herman, T., Fowler, T., Sheasha, T., Georgiou, T., Scialom, T., Speckbacher, T., Mihaylov, T., Xiao, T., Karn, U., Goswami, V., Gupta, V., Ramanathan, V., Kerkez, V., Gonguet, V., Do, V., Vogeti, V., Albiero, V., Petrovic, V., Chu, W., Xiong, W., Fu, W., Meers, W., Martinet, X., Wang, X., Wang, X., Tan, X. E., Xia, X., Xie, X., Jia, X., Wang, X., Goldschlag, Y., Gaur, Y., Babaei, Y., Wen, Y., Song, Y., Zhang, Y., Li, Y., Mao, Y., Coudert, Z. D., Yan, Z., Chen, Z., and Papakipos, Z. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- Kim, S., Hooper, C., Gholami, A., Dong, Z., Li, X., Shen, S., Mahoney, M. W., and Keutzer, K. Squeezellm: Dense-and-sparse quantization, 2024. URL <https://arxiv.org/abs/2306.07629>.
- Lin, J., Tang, J., Tang, H., Yang, S., Chen, W.-M., Wang, W.-C., Xiao, G., Dang, X., Gan, C., and Han, S. Awq: Activation-aware weight quantization for llm compression and acceleration, 2024. URL <https://arxiv.org/abs/2306.00978>.
- Merity, S., Xiong, C., Bradbury, J., and Socher, R. Pointer sentinel mixture models, 2016. URL <https://arxiv.org/abs/1609.07843>.
- Qwen, :, Yang, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Li, C., Liu, D., Huang, F., Wei, H., Lin, H., Yang, J., Tu, J., Zhang, J., Yang, J., Yang, J., Zhou, J., Lin, J., Dang, K., Lu, K., Bao, K., Yang, K., Yu, L., Li, M., Xue, M., Zhang, P., Zhu, Q., Men, R., Lin, R., Li, T., Tang, T., Xia, T., Ren, X., Ren, X., Fan, Y., Su,

Y., Zhang, Y., Wan, Y., Liu, Y., Cui, Z., Zhang, Z., and Qiu, Z. Qwen2.5 technical report, 2025. URL <https://arxiv.org/abs/2412.15115>.

Qwen Team. Qwen3.5. <https://huggingface.co/Qwen/Qwen3.5-9B>, 2026.

Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring the limits of transfer learning with a unified text-to-text transformer, 2023. URL <https://arxiv.org/abs/1910.10683>.

Shao, W., Chen, M., Zhang, Z., Xu, P., Zhao, L., Li, Z., Zhang, K., Gao, P., Qiao, Y., and Luo, P. Omniquant: Omnidirectionally calibrated quantization for large language models, 2024. URL <https://arxiv.org/abs/2308.13137>.

Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and Lample, G. Llama: Open and efficient foundation language models, 2023. URL <https://arxiv.org/abs/2302.13971>.

Tseng, A., Chee, J., Sun, Q., Kuleshov, V., and Sa, C. D. Quip#: Even better llm quantization with hadamard incoherence and lattice codebooks, 2024. URL <https://arxiv.org/abs/2402.04396>.

van Baalen, M., Kuzmin, A., Koryakovskiy, I., Nagel, M., Couperus, P., Bastoul, C., Mahurin, E., Blankevoort, T., and Whatmough, P. Gptvq: The blessing of dimensionality for llm quantization, 2025. URL <https://arxiv.org/abs/2402.15319>.

## A. $g=16$ Results on Qwen2.5

Table 6 presents  $g=16$  WikiText-2 results on the Qwen2.5 family. AAAC achieves the best result on all four models.

Table 6.  $g=16$  W4A16 perplexity on Qwen2.5 models, WikiText-2 ( $\downarrow$ ). Best quantized result in **bold**.

Model	BF16	RTN	IF4	AWQ	AAAC
Qwen2.5-0.5B	13.03	14.33	14.11	14.17	<b>13.66</b>
Qwen2.5-7B	6.80	7.03	6.97	7.00	<b>6.93</b>
Qwen2.5-14B	5.25	5.59	5.54	5.54	<b>5.45</b>
Qwen2.5-32B	4.97	5.17	5.12	5.15	<b>5.09</b>
<i>Average</i>	<i>7.51</i>	<i>8.03</i>	<i>7.94</i>	<i>7.97</i>	<i>7.78</i>
<i>Recovery</i>	—	—	18.4%	12.6%	<b>47.8%</b>

Table 7 presents the corresponding C4 results. The same pattern holds across all models.

Table 7.  $g=16$  W4A16 perplexity on Qwen2.5 models, C4 ( $\downarrow$ ). Best quantized result in **bold**.

Model	BF16	RTN	IF4	AWQ	AAAC
Qwen2.5-0.5B	20.02	21.82	21.57	21.54	<b>20.95</b>
Qwen2.5-7B	11.77	12.06	12.02	11.99	<b>11.93</b>
Qwen2.5-14B	10.25	10.52	10.49	10.45	<b>10.39</b>
Qwen2.5-32B	10.09	10.23	10.21	10.21	<b>10.18</b>
<i>Average</i>	<i>13.03</i>	<i>13.66</i>	<i>13.57</i>	<i>13.55</i>	<i>13.36</i>
<i>Recovery</i>	—	—	13.6%	17.6%	<b>47.2%</b>

## B. $g=16$ C4 Results

Table 8 presents C4 perplexity for the same models as Table 1. AAAC achieves the lowest perplexity on all models, with 50.9% average gap recovery.

Table 8.  $g=16$  W4A16 perplexity on C4 ( $\downarrow$ ). AAAC uses  $S=g=16$  (zero overhead). Best quantized result in **bold**.

Model	BF16	RTN	IF4	AWQ	AAAC
Llama-3.2-1B	13.79	15.60	15.30	15.15	<b>14.62</b>
Llama-3.2-3B	11.20	12.02	11.86	11.81	<b>11.58</b>
Llama-3.1-8B	9.50	10.21	10.07	10.07	<b>9.84</b>
Qwen3.5-2B	17.58	18.59	18.40	18.29	<b>18.13</b>
Qwen3.5-4B	13.99	14.51	14.37	14.38	<b>14.28</b>
Qwen3.5-9B	12.19	12.57	12.53	12.44	<b>12.38</b>
Qwen3.5-27B	10.47	10.64	10.61	10.59	<b>10.55</b>
<i>Average</i>	<i>12.67</i>	<i>13.45</i>	<i>13.31</i>	<i>13.25</i>	<i>13.05</i>
<i>Recovery</i>	—	—	18.5%	26.0%	<b>50.9%</b>

## C. Downstream Task Results

Table 9 reports per-task accuracy for Qwen3.5-9B and Qwen3.5-27B at  $g=16$ . All methods use the NVFP4 memory layout. Evaluations use the default few-shot settings from lm-evaluation-harness (Gao et al., 2024).

**AAAC: Activation-Aware Adaptive Codebooks for 4-bit LLM Quantization**

Table 9.  $g=16$  W4A16 downstream task accuracy ( $\% \uparrow$ ) on Qwen3.5-9B and Qwen3.5-27B. All methods use the NVFP4 memory layout. Best quantized average in **bold**.

Model	Task	BF16	RTN	IF4	AWQ	AAAC
Qwen3.5-9B	HellaSwag	78.1	77.3	77.3	77.3	77.6
	ARC-C	56.0	55.3	55.5	53.8	54.9
	PIQA	79.9	79.2	80.0	80.0	80.0
	BBH	62.3	60.6	60.1	61.4	61.7
	MuSR	43.4	42.1	40.5	44.2	44.2
	GPQA	42.2	39.7	41.9	41.4	40.0
	MATH	47.9	47.4	49.0	48.6	48.9
	MMLU	78.7	77.3	77.3	77.3	78.1
Qwen3.5-27B	HellaSwag	83.4	83.0	83.2	83.2	83.2
	ARC-C	61.7	61.9	61.3	61.9	61.9
	PIQA	82.1	82.3	82.2	81.8	82.5
	BBH	73.1	72.5	72.2	72.8	72.3
	MuSR	52.5	53.7	51.3	52.0	51.6
	GPQA	50.1	49.3	48.2	49.3	48.5
	MATH	57.6	55.7	57.4	56.1	57.4
	MMLU	84.5	84.2	84.5	84.2	84.7
<i>Average</i>		<i>64.59</i>	<i>63.84</i>	<i>63.87</i>	<i>64.08</i>	<b><i>64.21</i></b>
<i>Recovery</i>		—	—	4.0%	32.0%	<b>49.3%</b>

### D. W4A8 Results

Although AAAC uses activations to weight importance during calibration, it is not sensitive to activation precision at inference time. The same  $g=128$  quantized models from Table 2 are evaluated with per-tensor FP8 activation quantization. Table 10 confirms that AAAC’s improvements carry over to the W4A8 setting, with AAAC at  $S=16$  recovering 49.1% of the quantization gap.

Table 10.  $g=128$  W4A8 WikiText-2 perplexity ( $\downarrow$ ) vs. AWQ. Same quantized models as Table 2 with FP8 activations. **Bold**: best result per model.

Model	FP16	RTN	AWQ	AAAC	
				$S=128$	$S=16$
Llama-3.2-3B	7.85	8.54	8.26	<b>8.20</b>	<b>8.15</b>
Llama-3.1-8B	6.27	6.87	6.65	6.72	<b>6.60</b>
Qwen2.5-0.5B	13.22	15.74	15.05	<b>14.64</b>	<b>14.44</b>
Qwen2.5-7B	6.89	7.27	7.11	7.17	<b>7.11</b>
Qwen2.5-14B	5.33	5.82	5.70	<b>5.69</b>	<b>5.64</b>
<i>Average</i>	<i>7.91</i>	<i>8.85</i>	<i>8.55</i>	<i>8.48</i>	<b><i>8.39</i></b>
<i>Recovery</i>	—	—	31.4%	38.9%	<b>49.1%</b>