

GuidedQuant: Large Language Model Quantization via Exploiting End-Loss Guidance

Jinuk Kim^{1,2} Marwa El Halabi³ Wonyo Park⁴ Clemens JS Schaefer⁴
Deokjae Lee^{1,2} Yeonhong Park¹ Jae W. Lee¹ Hyun Oh Song^{1,2}

¹Seoul National University ²Neural Processing Research Center

³Samsung AI Lab, Montreal ⁴Google

ICML 2025

Demo: GuidedQuant

Llama-3.3-70B-Instruct (2-bit quantized)

Input: Write me a short and concise story about Harry, Ron, and Hermione.

SqueezeLLM

Harry, Ron, and Hermione were the best of friends.
They had been friends since they were kids, and their friendship had been strong since then.
They had been through a lot together, and their bond had grown even stronger.
They had been through a lot, and they had been through a lot of things, and they had been through a lot of things.
They had been through a lot, and they had been through a lot of things, and they had been through a lot of things.
They had been through a lot, and they had been through a lot of things, and



GuidedQuant + LNQ (Ours)

As the sun set over the Hogwarts grounds, Harry, Ron, and Hermione sat together on the grass, watching the stars twinkle to life.
They had just finished a long day of classes, and were now unwinding.
Ron began to tell them about his latest Quidditch match, and they all shared a laugh.
The three of them had been friends since their first year and had always been there to support each other.
Now, as they sat together, they each knew that they would always be there for one another.



Problem: Post-training Quantization for LLMs

There are two dominant PTQ approaches for LLM quantization:

- **Layer-wise Output-based Quantization:** Minimize the mean-squared error between each layer's original and quantized outputs.

$$\|\mathbf{X}^{(l)}\mathbf{W}^{(l)} - \mathbf{X}^{(l)}\widehat{\mathbf{W}}^{(l)}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^{d_{\text{out}}^{(l)}} \left(Z_{ij}^{(l)} - \widehat{Z}_{ij}^{(l)} \right)^2$$

Problem: Post-training Quantization for LLMs

There are two dominant PTQ approaches for LLM quantization:

- ▶ **Layer-wise Output-based Quantization:** Minimize the mean-squared error between each layer's original and quantized outputs.

$$\|\mathbf{X}^{(l)}\mathbf{W}^{(l)} - \mathbf{X}^{(l)}\widehat{\mathbf{W}}^{(l)}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^{d_{\text{out}}^{(l)}} \left(Z_{ij}^{(l)} - \widehat{Z}_{ij}^{(l)} \right)^2$$

- ▶ **Diagonal Fisher Information Matrix (FIM):** Weigh the individual weight errors by gradients of the end loss.

$$(\widehat{\mathbf{w}} - \mathbf{w})^\top \text{diag}(\mathbf{F})(\widehat{\mathbf{w}} - \mathbf{w}) = \sum_k \left(\frac{\partial \ell}{\partial w_k} \right)^2 (\widehat{w}_k - w_k)^2$$

Problem: Post-training Quantization for LLMs

- **Layer-wise Output-based Quantization:** Minimize the mean-squared error between each layer's original and quantized outputs.

$$\|\mathbf{X}^{(l)}\mathbf{W}^{(l)} - \mathbf{X}^{(l)}\widehat{\mathbf{W}}^{(l)}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^{d_{\text{out}}^{(l)}} \left(Z_{ij}^{(l)} - \widehat{Z}_{ij}^{(l)} \right)^2$$

→ Treats all hidden features equally, ignoring how much each feature contributes to the end loss.

Problem: Post-training Quantization for LLMs

- **Diagonal Fisher Information Matrix (FIM):** Weigh the individual weight errors by gradients of the end loss.

$$\begin{aligned}\ell(\hat{\mathbf{w}}) - \ell(\mathbf{w}) &\approx (\hat{\mathbf{w}} - \mathbf{w})^\top \mathbf{F}(\hat{\mathbf{w}} - \mathbf{w}) && \text{(2nd order Taylor expansion)} \\ &\approx (\hat{\mathbf{w}} - \mathbf{w})^\top \text{diag}(\mathbf{F})(\hat{\mathbf{w}} - \mathbf{w}) && \text{(Diagonal Approximation)} \\ &= \sum_k \left(\frac{\partial \ell}{\partial w_k} \right)^2 (\hat{w}_k - w_k)^2\end{aligned}$$

→ Approximates FIM with only its diagonal, discarding cross-weight interactions, which are crucial.

GuidedQuant: Objective

GuidedQuant bridges these gaps by **integrating end loss gradients** into the quantization objective while **preserving intra-channel dependencies**.

$$\left\| \frac{\partial \ell}{\partial \mathbf{Z}^{(l)}} \odot (\mathbf{X}^{(l)} \mathbf{W}^{(l)} - \mathbf{X}^{(l)} \widehat{\mathbf{W}}^{(l)}) \right\|_F^2 = n \sum_{l=1}^L \sum_{j=1}^{d_{\text{out}}^{(l)}} (\mathbf{w}_j^{(l)} - \widehat{\mathbf{w}}_j^{(l)})^\top \mathbf{F}_j^{(l)} (\mathbf{w}_j^{(l)} - \widehat{\mathbf{w}}_j^{(l)}).$$

$$\left\| \frac{\partial \ell}{\partial (\mathbf{X}\mathbf{W})} \in \mathbb{R}^{n \times d_{\text{out}}} \odot \left(\mathbf{X} \in \mathbb{R}^{n \times d_{\text{in}}} \left(\widehat{\mathbf{W}} \in \mathbb{R}^{d_{\text{in}} \times d_{\text{out}}} - \mathbf{W} \in \mathbb{R}^{d_{\text{in}} \times d_{\text{out}}} \right) \right) \right\|_F^2$$

$\underbrace{\hspace{1.5cm}}_{J_1} \quad \dots \quad \underbrace{\hspace{1.5cm}}_{J_g} \quad \quad \underbrace{\hspace{1.5cm}}_{J_1} \quad \dots \quad \underbrace{\hspace{1.5cm}}_{J_g}$

GuidedQuant: Objective

GuidedQuant bridges these gaps by **integrating end loss gradients** into the quantization objective while **preserving intra-channel dependencies**.

$$\left\| \frac{\partial \ell}{\partial \mathbf{Z}^{(l)}} \odot (\mathbf{X}^{(l)} \mathbf{W}^{(l)} - \mathbf{X}^{(l)} \widehat{\mathbf{W}}^{(l)}) \right\|_F^2 = n \sum_{l=1}^L \sum_{j=1}^{d_{\text{out}}^{(l)}} (\mathbf{w}_j^{(l)} - \widehat{\mathbf{w}}_j^{(l)})^\top \mathbf{F}_j^{(l)} (\mathbf{w}_j^{(l)} - \widehat{\mathbf{w}}_j^{(l)}).$$

This objective corresponds to **block-diagonal FIM approximation**, and is a more accurate approximation than

- **Layer-wise output error** objective which assumes $\frac{\partial \ell}{\partial \mathbf{Z}^{(l)}} \propto \mathbf{I}$.
- **Diagonal FIM** objective which ignores off-diagonal entries.

GuidedQuant: Averaging Approximation

- ▶ However, $\mathbf{H}_j^{(l)} := \mathbf{F}_j^{(l)}$ depends on each output channel $j \in \{1, \dots, d_{\text{out}}^{(l)}\}$, making it infeasible to compute and store.
- ▶ Solution: Group d_{out} channels into $g \ll d_{\text{out}}$ clusters (J_1, \dots, J_g) , and average Fisher blocks within each group.

$$\bar{\mathbf{H}}_k^{(l)} = \frac{1}{|J_k|} \sum_{j \in J_k} \mathbf{H}_j^{(l)}.$$

$$\forall j \in \{1, \dots, d_{\text{out}}\} : \quad \mathbf{H}_j = \mathbf{X}^\top \text{Diag}\left(\frac{\partial \ell}{\partial \mathbf{z}_j}\right)^2 \mathbf{X} \quad \left| \quad \forall k \in \{1, \dots, g\} : \quad (g \ll d_{\text{out}}) \quad \bar{\mathbf{H}}_k = \frac{1}{|J_k|} \sum_{j \in J_k} \mathbf{H}_j\right.$$

The diagram illustrates the averaging approximation for Fisher blocks. It is divided into two parts by a vertical dashed line.

Left part (individual channel): For $\forall j \in \{1, \dots, d_{\text{out}}\}$, the Fisher block \mathbf{H}_j is computed as $\mathbf{H}_j = \mathbf{X}^\top \text{Diag}\left(\frac{\partial \ell}{\partial \mathbf{z}_j}\right)^2 \mathbf{X}$. This is visualized as a 3x2 grid of blue blocks (representing $\widehat{\mathbf{W}}_{:,j}$ and $\mathbf{W}_{:,j}$) followed by a 3x3 grid of red blocks (representing \mathbf{H}_j).

Right part (grouped channels): For $\forall k \in \{1, \dots, g\}$ where $g \ll d_{\text{out}}$, the average Fisher block $\bar{\mathbf{H}}_k$ is computed as $\bar{\mathbf{H}}_k = \frac{1}{|J_k|} \sum_{j \in J_k} \mathbf{H}_j$. This is visualized as a 3x2 grid of blue blocks (representing $\widehat{\mathbf{W}}_{:,J_k}$ and $\mathbf{W}_{:,J_k}$) followed by a 3x3 grid of green blocks (representing $\bar{\mathbf{H}}_k$).

GuidedQuant: Results

GuidedQuant can be plugged into any layer-wise PTQ backend.

	Method	Bits↓	Wiki2-4K↓
Type	Original	16	5.12
Weight-only Scalar	SqueezeLLM	2.01	39.58
	LNQ (Ours)	2.01	23.31
	LNQ + GQuant (Ours)	2.01	8.83
Weight-only Vector	QTIP	2.00	6.82
	QTIP + GQuant (Ours)	2.00	6.11
	Method	Bits↓	Wiki2-2K↓
Type	Original	16	5.47
Weight-and- Activation	SpinQuant	W4A4KV4	5.95
	SpinQuant + GQuant (Ours)	W4A4KV4	5.89

→ Improves state-of-the-art methods for **weight-only scalar**, **weight-only vector**, and **weight-and-activation quantization**.

LNQ: Problem Formulation

Regarding **weight-only scalar quantization**, we further propose non-uniform quantization method **LNQ (Layer-wise Non-uniform Quantization)**.

The optimization problem involves discrete assignment $\mathbf{P}^{(j)} \in \{0, 1\}^{d_{\text{in}} \times m}$ and continuous codebook $\mathbf{c}^{(j)} \in \mathbb{R}^m$:

$$\begin{aligned} & \underset{\substack{\mathbf{P}^{(j)} \in \{0,1\}^{d_{\text{in}} \times m} \\ \mathbf{c}^{(j)} \in \mathbb{R}^m}}{\text{minimize}} && \sum_{j=1}^{d_{\text{out}}} \|\mathbf{X}\mathbf{w}_j - \mathbf{X}\mathbf{P}^{(j)}\mathbf{c}^{(j)}\|_2^2 \\ & \text{subject to} && \mathbf{P}^{(j)}\mathbf{1}_m = \mathbf{1}_{d_{\text{in}}}, \end{aligned}$$

LNQ: Algorithm

LNQ is a non-uniform scalar quantization method that alternates **closed-form codebook update** and **coordinate-descent assignment update**.

repeat:

$$\mathbf{c}^{(j)} \leftarrow \left(\mathbf{P}^{(j)\top} \mathbf{X}^\top \mathbf{X} \mathbf{P}^{(j)} \right)^{-1} \mathbf{P}^{(j)\top} \mathbf{X}^\top \mathbf{X} \mathbf{w}_j \quad (\text{codebook})$$

for $i = 1$ **to** d_{in} **do:**

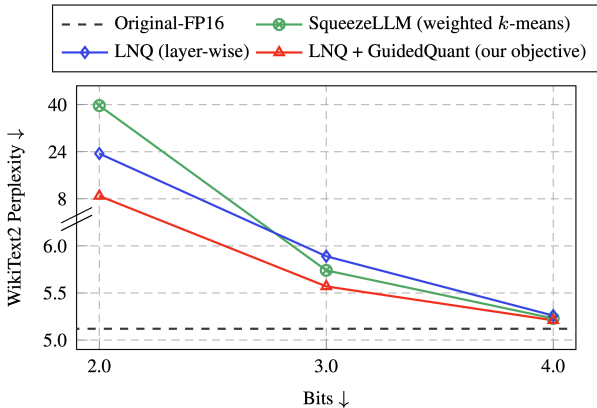
$$c_{q^*}^{(j)} \leftarrow \underset{\widehat{W}_{ij} \in \{c_1^{(j)}, \dots, c_m^{(j)}\}}{\operatorname{argmin}} (\widehat{\mathbf{w}}_j - \mathbf{w}_j)^\top \mathbf{X}^\top \mathbf{X} (\widehat{\mathbf{w}}_j - \mathbf{w}_j)$$

$$\forall q \in [m] : P_{iq}^{(j)} = \begin{cases} 1 & \text{if } q = q^*, \\ 0 & \text{otherwise.} \end{cases} \quad (\text{assignment})$$

→ Monotonically decreases the objective and guarantees convergence.

LNQ: Results

LNQ algorithm is fully compatible with **GuidedQuant**: Together, they achieve state-of-the-art performance on **weight-only scalar quantization**.



Conclusion

- ▶ **GuidedQuant** integrates end loss gradients into the layer-wise quantization objective, outperforming PTQ methods.
- ▶ **LNQ** is a non-uniform scalar quantization method that alternates closed-form codebook update and coordinate-descent assignment update.



Code: <https://github.com/snu-mllab/GuidedQuant>