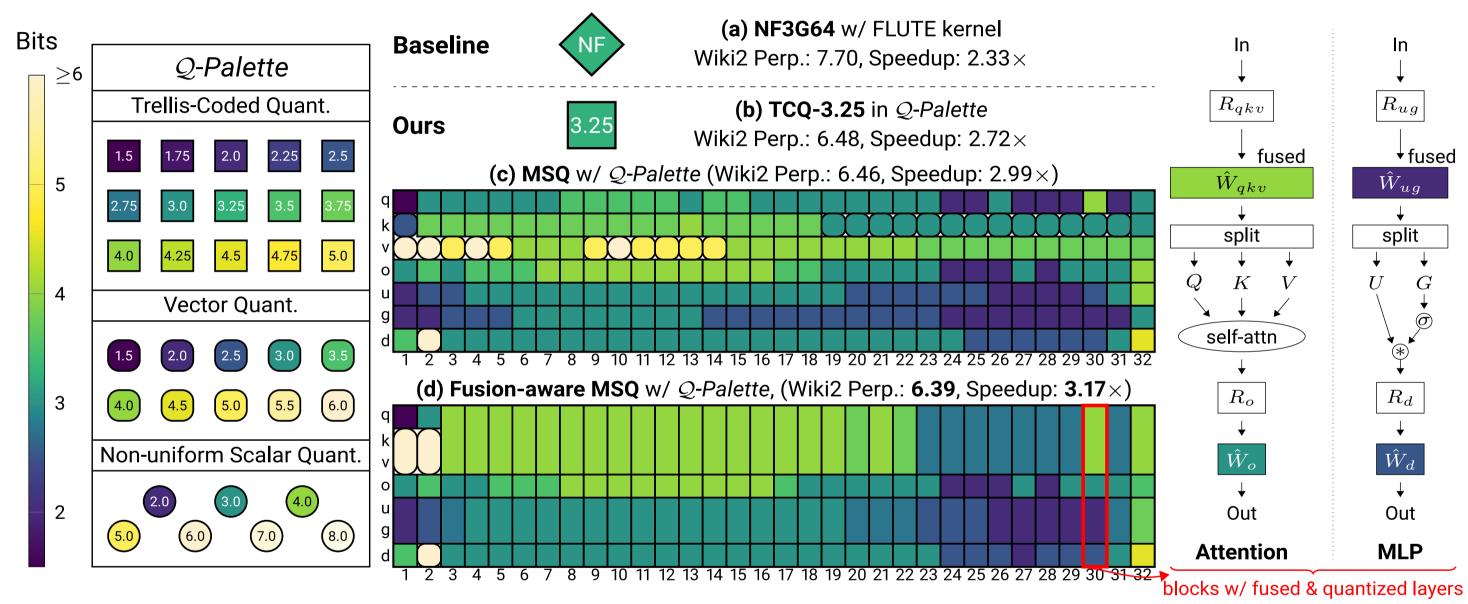
# Q-Palette: Fractional-Bit Quantizers Toward Optimal Bit Allocation for Efficient LLM Deployment

Deokjae Lee, Hyun Oh Song

**Seoul National University** 



### **Overview**



# **Contributions**

- 1. We develop **Q-Palette**, a set of versatile rotation-based quantizers with efficient inference CUDA kernels and wide fractional-bit support.
- 2. Built on Q-Palette, we propose a novel mixed-scheme quantization framework that **jointly optimizes quantizer selection and layer fusion**.

### **Results**

# More results are in our paper! | Samuration | Samuration

Figure 1. Performance trade-offs under different constraints (LLaMA 3.1-8B, RTX4090 GPU).

### **Motivation**

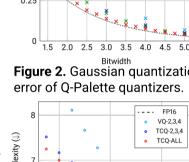
- Recent quantization methods utilize rotation to Gaussianize weights, reducing outliers.
- Based on HIGGS's linearized surrogate loss

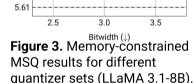
$$\mathcal{L}(\{Q(W_l)\}_{l=1}^L) - \mathcal{L}(\{W_l\}_{l=1}^L) pprox \sum_{l=1}^L a_l \underbrace{\|Q(W_l) - W_l\|^2 / \|W_l\|^2}_{=: \, ext{err}(Q; W_l)},$$

the memory-constrained mixed scheme quantization (MSQ) can be formally written by

$$egin{align*} & \min_{P_{lq} \in \{0,1\}} & \sum_{l=1}^{L} \sum_{q=1}^{|\mathcal{Q}|} P_{lq} \cdot \underbrace{\left(a_l \cdot \operatorname{err}(Q_q; W_l)
ight)}_{ ext{loss term $\ell_{lq}$}} \ & \operatorname{subject\ to} & \sum_{q=1}^{|\mathcal{Q}|} P_{lq} = 1 \quad orall 1 \leq l \leq L, \ & \sum_{l=1}^{L} \sum_{q=1}^{|\mathcal{Q}|} P_{lq} \cdot \operatorname{bit}(Q_q; W_l) d_l^{\operatorname{in}} d_l^{\operatorname{out}} \leq M. \end{cases}$$

- Here, for Gaussianized weights, we have a natural lower bound on distortion term  $\mathbb{E}[\text{err}(Q)] \geq 2^{-2\text{bit}(Q)}$ , given by Gaussian distortion-rate theorem.
- Assuming ideal Gaussian quantizers that achieve the optimal distortion bound, the optimal bit allocation is given by





$$b_l^* = \max iggl\{ 0, rac{1}{2\ln(2)} iggl( \ln rac{a_l}{d_l^{ ext{in}} d_l^{ ext{out}}} iggr) + C iggr\} \quad orall \, 1 \leq l \leq L, \quad ext{for } \, C \, ext{ s.t. } \, \sum_{l=1}^L b_l^* d_l^{ ext{in}} d_l^{ ext{out}} = M.$$

However, in practice, quantization is performed with a set of non-ideal quantizers, and the
performance gap is given by 1) how closely each quantizer approaches the distortion bound,
and 2) how finely the available bitwidths approximate the optimal bit allocations, motivating
the design of Q-Palette (Sec 3.2).

## **Fusion-aware MSQ**

- We propose **fusion-aware MSQ**, a novel MSQ framework that jointly optimizes quantization with the additional design dimension of **layer fusion**. (see Overview)
- Fusion-aware MSQ simultaneously determines 1) how to group layers for fusion and 2) which quantizer to assign to each fused group, encoded by  $P_{gg}$ .
- The fusion-aware MSQ problem is formulated as

$$\begin{array}{ll} \underset{P_{gq} \in \{0,1\}}{\operatorname{minimize}} & \sum_{g \in \mathcal{G}} \sum_{q=1}^{|\mathcal{Q}|} P_{gq} \cdot \sum_{l \in g} \ell_{lq} \\ \text{subject to} & \sum_{g \in \mathcal{G}: l \in g} \sum_{q=1}^{|\mathcal{Q}|} P_{gq} = 1 \quad \forall 1 \leq l \leq L, \quad \sum_{g \in \mathcal{G}} \sum_{q=1}^{|\mathcal{Q}|} P_{gq} \cdot c_{gq} \leq C \end{array}$$

- Here,  $\mathcal G$  denotes the set of all fusible layer groups, where each group consists of linear layers sharing the same input (e.g., Q, K, V layers in the same Transformer block).
- The problem above is ILP and can be solved by solvers such as SCIP solver.