

# RESQ : MIXED PRECISION QUANTIZATION OF LARGE LANGUAGE MODELS WITH LOW-RANK RESIDUALS

Utkarsh Saxena<sup>1\*</sup>, Sayeh Sharify<sup>2</sup>, Kaushik Roy<sup>1</sup>, Xin Wang<sup>2</sup>

<sup>1</sup>Purdue University, West Lafayette, USA

<sup>2</sup>d-Matrix, Santa Clara, USA

## ABSTRACT

Quantizing weights, activations, and KV cache in large language models to 4-bit without degrading generalizability is challenging due to outlier-induced activation quantization errors. We propose *ResQ*, a post training quantization (PTQ) method that uses principal component analysis to identify a low-rank subspace (in practice  $1/8$  of the hidden dimension) and keeps coefficients within this subspace in 8-bit while quantizing the rest in 4-bit. Within each subspace, invariant random rotation is applied to further suppress outliers. ResQ outperforms recent PTQ methods on Llama and Qwen2.5, achieving up to 33% lower Wikitext perplexity than *SpinQuant* and up to  $3\times$  speedup over 16-bit. Code is available at here<sup>1</sup>.

## 1 INTRODUCTION

Quantization enables efficient on-device large language model (LLM) inference by reducing storage (weight quantization), memory usage (KV cache quantization), and compute complexity (activation quantization). While post-training methods achieve 2-bit KV cache quantization (Liu et al., 2024b) and low-precision weights (Frantar et al., 2022), activation quantization below 8-bit remains challenging due to activation outliers (Dettmers et al., 2022). Recent methods employ two key strategies: (1) *Differential treatment of outliers*, where select channels are preserved in high precision, yielding mixed-precision quantization (Dettmers et al., 2022; Zhao et al., 2024; Ashkboos et al., 2024b). QUIK (Ashkboos et al., 2024b) and ATOM (Zhao et al., 2024) statically retain outlier channels in 8-bit. (2) *Invariant random rotation*, which suppresses activation outliers for uniform low-precision quantization (Ashkboos et al., 2024c; Liu et al., 2024a). QuaRot (Ashkboos et al., 2024c) applies Hadamard rotations to activations, while SpinQuant (Liu et al., 2024a) optimizes activation rotations via gradient descent. Both types of activation quantization approaches reduce quantization error; yet a notable model performance gap persists from the 16-bit baseline.

To address this gap, we introduce *ResQ*, a novel PTQ method for efficient 4-bit quantization of activations, weights, and KV cache. Using offline principal component analysis (PCA), ResQ identifies a low-rank subspace capturing highest variance in activations, quantizing its coefficients in 8-bit while applying 4-bit quantization to the rest. ResQ then employs invariant random rotations within each subspace to further suppress outliers, minimizing error with most projections fused into adjacent weights for minimal overhead. Compared with related activation quantization approaches, ResQ achieves highest quantization SNR (Figure 1(b)) with its provably optimal choice of components in 8-bit. It supports KV cache quantization and integrates with GPTQ (Frantar et al., 2022), enhancing LLM generalization. With only  $1/8$  channels in 8-bit, ResQ reduces perplexity by 4–33% on Wikitext and improves 0-shot accuracy by 0.1–5.4% over SpinQuant (Liu et al., 2024a), without requiring gradient-based optimization. Compared with 16-bit floating point model, ResQ achieves upto  $3\times$  inference speedup which is only on an average 14% less than fully INT4 inference.

We claim the following contributions : (1) We propose ResQ, a mixed precision weight, activation, and KV cache quantization method by keeping low-rank, high-variance components in high precision, in combination with random rotation-induced outlier suppression. (2) We theoretically analyze the projection matrices in ResQ and show that using PCA-based projections minimizes quantization

\*corresponding author : Utkarsh Saxena (saxenau@purdue.edu)

<sup>1</sup><https://github.com/utkarsh-dmx/project-resq>

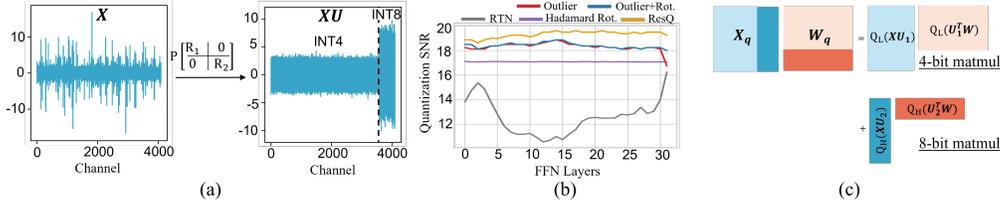


Figure 1: (a) Activation distribution before and after ResQ projections, (b) Quantization SNR for Feed forward network layers (higher is better), (c) Mixed precision matrix multiplication.

error. (3) We conduct extensive experiments on various models and language tasks and show that ResQ outperforms related state-of-the-art approaches. (4) We develop CUDA kernels and achieve runtime speedup on NVIDIA GPUs with our quantized models.

## 2 RESQ

In this section, we present ResQ, a mixed-precision quantization method that preserves low-rank components in 8-bit while quantizing the rest at lower precision, with theoretical guarantees and efficient LLM deployment.

### 2.1 QUANTIZATION SCHEME AND PROJECTIONS

Given input activation  $X \in \mathbb{R}^{n \times d}$  and weight  $W \in \mathbb{R}^{d \times d}$ , we project them onto orthogonal basis  $U \in \mathbb{R}^{d \times d}$ , then quantize the coefficients. High-precision components are captured by  $U_h \in \mathbb{R}^{d \times r}$ , and low precision by  $U_l \in \mathbb{R}^{d \times (d-r)}$ , ensuring  $U_h U_h^T + U_l U_l^T = U U^T = I$ . The rank  $r$  controls the amount of components in high precision (in practice we typically choose  $r = d/s$ ). Where  $Q(\cdot)$  is the quantization operator, the quantized activation and weights are,

$$X_q = Q_L(XU_l) + Q_H(XU_h), \quad W_q = Q_L(U_l^T W) + Q_H(U_h^T W) \quad (1)$$

The layer output is given below and also demonstrated in Figure 1(c).

$$X_q W_q = Q_L(XU_l)Q_L(U_l^T W) + Q_H(XU_h)Q_H(U_h^T W). \quad (2)$$

Due to orthogonality, the projections preserve the original model output in absence of quantization. The orthogonal basis  $U$  should (1) prioritize important components for high-precision quantization and (2) minimize quantization error in both high- and low-precision groups. We construct  $U = PR = [P_l P_h] \begin{bmatrix} R_l & 0 \\ 0 & R_h \end{bmatrix}$  using two rotation matrices:  $P$  for importance based projections and  $R$  to minimize quantization error. Inspired by prior work (Ashkboos et al., 2024c; Chee et al., 2024), we make  $R_l, R_h$  random orthogonal matrices because random rotation reduces outliers, making the rotated matrices easier to quantize. Furthermore, projection with a random orthogonal matrix increases Gaussianity of activations and weights Tseng et al. (2024) within high- and low-precision groups, conducive to the quantizations applied to these groups. To determine  $P$ , we minimize the activation quantization error  $\|X - X_q\|_F$ . For activations quantized according to Equation 1, we have,

$$\|X - X_q\|_F = \|XU_l - Q_L(XU_l)\|_F + \|XU_h - Q_H(XU_h)\|_F. \quad (3)$$

**Theorem 2.1.** For any matrix  $X$  quantized to  $X_q$  according to method described in Equation 1, assuming the values to be quantized in  $X$  are normally distributed, we have

$$\mathbb{E}\|X - X_q\|_F \leq \frac{\sqrt{\pi \log(d-r)}}{2^{L-1} - 1} \mathbb{E}\|X\|_F - \left[ \frac{\sqrt{\pi \log(d-r)}}{2^{L-1} - 1} - \frac{\sqrt{\pi \log r}}{2^{H-1} - 1} \right] \mathbb{E}\|X P_h\|_F. \quad (4)$$

Full proof of Theorem 2.1 is in Appendix B. Theorem 2.1 bounds the quantization error, which can be minimized by maximizing  $\|X P_h\|_F$ . This occurs when  $P_h$  consists of top eigenvectors of covariance matrix of activations  $XX^T$ . Thus, the low rank space for high precision quantization can be obtained by means of PCA. To facilitate that, we obtained  $XX^T$  using a calibration dataset and perform PCA to obtain  $P_h$ . Note that this is done offline and once obtained,  $P_h$  does not change. The subspace for low-precision quantization can be obtained using

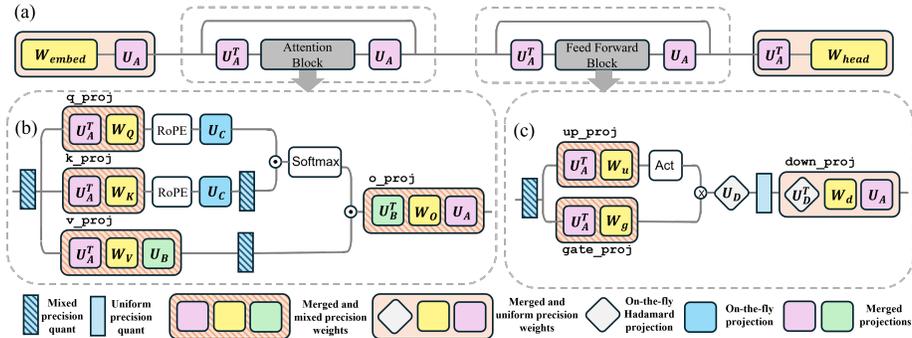


Figure 2: Model inference with ResQ incorporating the projection matrices. (a)  $U_A$  modifies the inputs across blocks enabling better quantization. (b)  $U_B, U_C$  enables mixed precision quantization of KV cache. (c)  $U_D$  projects the activations and weights of down\_proj layer.

$U_h U_h^\top + U_l U_l^\top = P_h P_h^\top + P_l P_l^\top = I$  (because  $R_i$  is orthogonal). If we construct  $P$  by taking eigenvectors of  $XX^\top$  arranged in *increasing* order of eigenvalues, the last  $r$  columns of such a  $P$  would correspond to  $P_h$  and the first  $d - r$  columns would correspond to  $P_l$ . The distribution of activation after before and after projection is given in Figure 1(a).

2.2 INFERENCE COMPUTATION WITH OPTIMIZED PROJECTIONS

Once projection matrices are obtained, activations and weights are projected using  $U$ . Weights are projected and quantized offline, while activation projections are merged into the previous layer’s weights for efficiency. Decoder-based LLMs require four projections (Figure 2):  $U_A$  (hidden dim),  $U_B, U_C$  (attention head dim), and  $U_D$  (FFN hidden dim). **Block Boundary Projections** Activations for attention and feed forward network block are projected via  $U_A$  by right-multiplying the final linear layer weights (o\_proj in attention, down\_proj in FFN), incurring no extra inference cost. To ensure numerical invariance, the first linear layers (q\_proj | k\_proj | v\_proj in attention, up\_proj | gate\_proj in FFN) are pre-multiplied with  $U_A^\top$ . Embedding and final head weights are also adjusted for residual stream projection. **Attention Block Projections**  $U_B, U_C$  project activations within the attention block (Figure 2b). Post-multiplying the value projection layer by  $U_B$  ensures optimal KV cache quantization, requiring o\_proj weights to be pre-multiplied by  $U_B^\top$  for numerical invariance.  $U_C$  optimally quantizes keys by projecting both query and key, preserving the attention dot product:  $q_{proj} K_{proj}^\top = (q U_C)(U_C^\top K^\top) = q K^\top$ . Since  $U_C$  cannot be merged due to RoPE, the projection is explicitly computed at runtime, but made more efficient by applying uniform precision quantization to  $U_C$  and corresponding inputs. **FFN Block Projections**  $U_D$  improves FFN activation quantization (Figure 2c). While  $U_D^\top$  is fused with down\_proj weights, activation functions prevent merging  $U_D$  with preceding layers, requiring runtime computation. Given the large FFN hidden dimension ( $d_{FFN}$ ), direct multiplication is costly. To mitigate this,  $U_D$  is a Hadamard matrix for efficient transforms. And down\_proj weights and activations are kept entirely in 4-bit.

3 EXPERIMENTS

**Models, Tasks, Datasets, and Baselines** We evaluate ResQ on Llama 3 (Meta, 2024b), Llama 3.2 (Meta, 2024a), Qwen2.5 (Yang et al., 2024a), and multi-modal Qwen2 VL models (Wang et al., 2024). Baselines include GPTQ (Frantar et al., 2022), QuaRot (Ashkboos et al., 2024c), QUIK (Ashkboos et al., 2024b), SpinQuant (Liu et al., 2024a), and SmoothQuant+, a stronger baseline created by combining SmoothQuant (Xiao et al., 2023) with GPTQ following Sharify et al. (2024). We assess quantization on *language modeling*

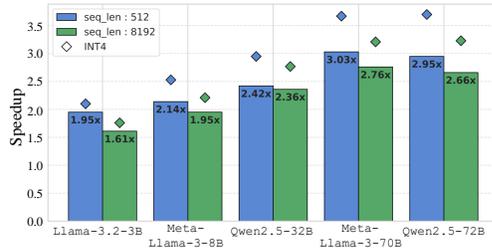


Figure 3: Decoder block speedup of ResQ and INT4 kernel on NVIDIA RTX 3090 over 16-bit floating point baseline for batch size of 1.

Table 1: Perplexity on Wikitext, average 0-shot common sense reasoning accuracy, and average 0-shot MMLU accuracy at W/A/KV = 4-bit. \*: Mixed precision with 1/8 channels in 8-bit for W/A/KV = 4.5-bit. †: higher is better, ‡: lower is better. Full results in Appendix D.

Llama 3						
Method	Meta-Llama-3-8B			Meta-Llama-3-70B		
	Wiki (‡)	Reasoning (†)	MMLU (†)	Wiki (‡)	Reasoning (†)	MMLU (†)
16-bit	6.1	67.1	63.1	2.9	73.1	75.9
RTN	218.9	39.3	23.6	452.7	45.5	23.2
GPTQ	166.3	39.8	23.3	1.2e4	34.9	25.5
SQ+	78.2	42.5	24.7	-	-	-
QUIK*	14.2	51.6	32.7	8.0	58.2	51.1
QuaRot	7.8	62.1	53.2	5.7	67.6	65.3
SpinQuant	7.4	63.8	56.2	6.2	65.7	59.4
ResQ*	<b>7.1</b>	<b>63.9</b>	<b>57.2</b>	<b>4.1</b>	<b>71.1</b>	<b>73.9</b>

Llama 3.2						
Method	Llama-3.2-1B			Llama-3.2-3B		
	Wiki (‡)	Reasoning (†)	MMLU (†)	Wiki (‡)	Reasoning (†)	MMLU (†)
16-bit	9.8	54.9	36.9	7.8	62.7	54.8
RTN	329.1	38.1	23.8	268.8	38.7	25.7
GPTQ	108.9	38.0	24.9	178.3	40.3	24.8
SQ+	228.9	38.0	24.1	96.1	39.0	25.9
QUIK*	21.8	44.3	25.1	15.8	48.8	31.1
QuaRot	14.3	49.0	25.5	10.1	56.1	42.0
SpinQuant	13.6	48.8	25.6	9.2	57.9	44.2
ResQ*	<b>12.4</b>	<b>50.1</b>	<b>29.4</b>	<b>8.8</b>	<b>59.0</b>	<b>49.8</b>

Qwen2.5						
Method	Qwen2.5-3B			Qwen2.5-72B		
	Wiki (‡)	Reasoning (†)	MMLU (†)	Wiki (‡)	Reasoning (†)	MMLU (†)
16-bit	8.0	63.8	66.1	3.9	73.4	84.3
RTN	5.9e4	35.1	23.4	4.5e4	34.3	24.0
GPTQ	9.9e3	35.1	23.2	3.8e4	34.5	23.3
SQ+	7.3e4	34.8	23.9	-	-	-
QUIK*	15.5	51.2	39.4	8.3	61.9	69.3
QuaRot	68.8	47.7	28.9	4.9	70.3	80.1
ResQ*	<b>9.0</b>	<b>61.1</b>	<b>61.2</b>	<b>4.6</b>	<b>72.0</b>	<b>81.5</b>

Table 2: Comparison of performance of quantization approaches on generative tasks at precisions of W/A/KV = 4-bit. \*: Mixed precision with 1/8 channels in 8-bit for W/A/KV = 4.5-bit. †: higher is better.

Model	Method	GSM8K 5-shot (†)		LongBench (†)		
		flexible-c	strict-m	qmsum	samsun	repo-p
Meta-Llama-3-8B	16-bit	51.0	50.6	23.9	44.8	66.4
	QUIK*	2.3	0.0	10.5	25.2	37.6
	QuaRot	27.6	27.1	22.0	43.8	60.6
	SpinQuant	29.8	29.6	23.0	43.9	<b>62.6</b>
	ResQ*	<b>33.6</b>	<b>33.2</b>	<b>23.1</b>	<b>44.1</b>	<b>62.3</b>
Llama-3.2-8B	16-bit	25.1	24.9	23.1	43.0	64.4
	QUIK*	2.5	0.0	15.9	31.7	30.9
	QuaRot	10.1	9.1	20.6	39.5	56.8
	SpinQuant	11.6	11.4	<b>21.7</b>	41.9	59.1
	ResQ*	<b>17.1</b>	<b>16.7</b>	<b>21.7</b>	<b>43.0</b>	<b>61.5</b>

Table 3: 0-shot MMMU accuracy (higher is better) of vision language models when quantized using various approaches. \*: Mixed precision with 1/8 channels in 8-bit and rest in 4-bit.

W/A/KV	Method	Model	
		Qwen2-VL-2B-Instruct	Qwen2-VL-7B-Instruct
16/16/16	Baseline	39.6	51.6
	RTN	25.0	26.7
4/4/4	GPTQ	27.7	24.9
	QuaRot	24.0	24.5
	QUIK*	26.3	28.9
	ResQ*	<b>29.7</b>	<b>47.0</b>
4.5/4.5/4.5	RTN	24.9	25.2
	GPTQ	23.4	24.3
	QuaRot	26.5	24.5
	QUIK*	28.4	26.4
	ResQ*	<b>34.0</b>	<b>48.8</b>

(Wikitext (Merity et al., 2016)), *reasoning* (average 0-shot accuracy on Arc-c/e (Clark et al., 2018), BoolQ (Clark et al., 2019), HellaSwag (Zellers et al., 2019), Openbook QA (Mihaylov et al., 2018), PIQA (Bisk et al., 2020), SIQA (Sap et al., 2019), WinoGrande (Sakaguchi et al., 2021)), *understanding* (MMLU (Hendrycks et al., 2021)), *math* (GSM8K (Cobbe et al., 2021)), *summarization* (samsun, qmsun from LongBench (Bai et al., 2024)), *code completion* (repobench-p (Liu et al., 2023b)), and *multi-modal* (MMMU (Yue et al., 2024)). More implementation details in App. C.

**Language modeling, understanding, and reasoning tasks** The results are presented in Table 1. We find that ResQ closes the gap to 16-bit performance and surpasses all quantization baselines across tasks and models. On Llama 3/3.2, ResQ achieves 4–33% lower Wikitext perplexity, 0.1–5.4% higher average 0-shot accuracy, and 1–14.5% better MMLU accuracy than SpinQuant, without additional training. For Qwen-2.5, ResQ outperforms all baselines, which fail to achieve competitive results. Compared to QUIK, another mixed precision approach, ResQ improves Wikitext perplexity by 42–50%, 0-shot accuracy by 5.8–12.3%, and MMLU accuracy by 4.3–24.5%. Full results are in Appendix D. **Generative tasks** We evaluate ResQ on auto-regressive tasks, including GSM8K (math), dialogue summarization (qmsun, samsun), and code completion (repobench-p, Table 2), to assess generation across domains. On GSM8K, where QUIK fails, ResQ surpasses SpinQuant by 3.8% (8B) and 5.5% (3B), narrowing the gap to 16-bit. In LongBench tasks, ResQ outperforms SpinQuant without additional training. **Multi-modal understanding** We evaluate ResQ on vision-language models (VLMs) by quantizing the Qwen2 VL family and testing on MMMU (Table 3, Yue et al. 2024). Only the language model is quantized, as it contains most parameters (over 10× for Qwen2-VL-7B-Instruct). ResQ outperforms baselines on 2B and 7B models, demonstrating superior accuracy and generalizability. Individual MMMU task results are in Appendix E. **Hardware Performance** We implement mixed-precision quantization using CUDA 11.8 and PyTorch, leveraging CUTLASS Thakkar et al. (2023) for INT4/INT8 GEMM on TensorCore. On an NVIDIA RTX 3090, ResQ achieves 1.61×–3.03× speedup over 16-bit for a single decoder block (Figure 3), with greater gains on larger models and shorter sequences. ResQ is only 14% slower than INT4, demonstrating minimal overhead from mixed precision and on-the-fly projections.

## 4 CONCLUSION

We introduce *ResQ*, a mixed-precision PTQ method for 4-bit LLM quantization. ResQ projects weights, activations, and KV cache to subspaces spanned by principal components, preserving high-variance components ( $\frac{1}{8}$  of hidden dimension) in 8-bit and quantizing the rest to 4-bit. It outperforms both uniform and mixed-precision baselines, demonstrating effectiveness across diverse tasks on Llama and Qwen models. Compared to SpinQuant, ResQ reduces WikiText perplexity by up to 33% without retraining and achieves up to  $3\times$  speedup over 16-bit.

## REFERENCES

- Saleh Ashkboos, Maximilian L Croci, Marcelo Gennari do Nascimento, Torsten Hoefler, and James Hensman. SliceGPT: Compress large language models by deleting rows and columns. *arXiv:2401.15024*, 2024a.
- Saleh Ashkboos, Iliia Markov, Elias Frantar, Tingxuan Zhong, Xincheng Wang, Jie Ren, Torsten Hoefler, and Dan Alistarh. QUIK: Towards end-to-end 4-bit inference on generative large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 3355–3371. Association for Computational Linguistics, 2024b. doi: 10.18653/v1/2024.emnlp-main.197. URL <https://aclanthology.org/2024.emnlp-main.197/>.
- Saleh Ashkboos, Amirkeivan Mohtashami, Maximilian Croci, Bo Li, Pashmina Cameron, Martin Jaggi, Dan Alistarh, Torsten Hoefler, and James Hensman. Quarot: Outlier-free 4-bit inference in rotated llms. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 100213–100240. Curran Associates, Inc., 2024c. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/file/b5b939436789f76f08b9d0da5e81af7c-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/b5b939436789f76f08b9d0da5e81af7c-Paper-Conference.pdf).
- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. LongBench: A bilingual, multitask benchmark for long context understanding. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3119–3137. Association for Computational Linguistics, 2024. doi: 10.18653/v1/2024.acl-long.172. URL <https://aclanthology.org/2024.acl-long.172>.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. PIQA: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7432–7439, 2020.
- Chi-Chih Chang, Wei-Cheng Lin, Chien-Yu Lin, Chong-Yan Chen, Yu-Fang Hu, Pei-Shuo Wang, Ning-Chi Huang, Luis Ceze, Mohamed S Abdelfattah, and Kai-Chiang Wu. Palu: Compressing kv-cache with low-rank projection. *arXiv:2407.21118*, 2024.
- Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher M De Sa. QuIP: 2-bit quantization of large language models with guarantees. *Advances in Neural Information Processing Systems*, 36, 2024.
- Jungwook Choi, Zhuo Wang, Swagath Venkataramani, Pierce I-Jen Chuang, Vijayalakshmi Srinivasan, and Kailash Gopalakrishnan. PACT: Parameterized clipping activation for quantized neural networks. *arXiv:1805.06085*, 2018.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv:1905.10044*, 2019. URL <https://arxiv.org/abs/1905.10044>.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457*, 2018.

- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv:2110.14168*, 2021.
- Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Gpt3. int8 (): 8-bit matrix multiplication for transformers at scale. *Advances in Neural Information Processing Systems*, 35: 30318–30332, 2022.
- Tim Dettmers, Ruslan Svirschevski, Vage Egiazarian, Denis Kuznedelev, Elias Frantar, Saleh Ashkboos, Alexander Borzunov, Torsten Hoefler, and Dan Alistarh. SpQR: A sparse-quantized representation for near-lossless llm weight compression. *arXiv:2306.03078*, 2023.
- Shichen Dong, Wen Cheng, Jiayu Qin, and Wei Wang. QAQ: Quality adaptive quantization for llm kv cache. *arXiv:2403.04643*, 2024.
- Vage Egiazarian, Andrei Panferov, Denis Kuznedelev, Elias Frantar, Artem Babenko, and Dan Alistarh. Extreme compression of large language models via additive quantization. *arXiv:2401.06118*, 2024.
- Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. GPTQ: Accurate post-training quantization for generative pre-trained transformers. *arXiv:2210.17323*, 2022.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muenighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, 07 2024. URL <https://zenodo.org/records/12608602>.
- Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference. In *Low-Power Computer Vision*, pp. 291–326. Chapman and Hall/CRC, 2022.
- Ziyi Guan, Hantao Huang, Yupeng Su, Hong Huang, Ngai Wong, and Hao Yu. APTQ: Attention-aware post-training mixed-precision quantization for large language models. In *Proceedings of the 61st ACM/IEEE Design Automation Conference*, pp. 1–6, 2024.
- Yefei He, Luoming Zhang, Weijia Wu, Jing Liu, Hong Zhou, and Bohan Zhuang. ZipCache: Accurate and efficient kv cache quantization with salient token identification. *arXiv:2405.14256*, 2024.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- Coleman Hooper, Sehoon Kim, Hiva Mohammadzadeh, Michael W Mahoney, Yakun Sophia Shao, Kurt Keutzer, and Amir Gholami. KVQuant: Towards 10 million context length llm inference with kv cache quantization. *arXiv:2401.18079*, 2024.
- Wei Huang, Haotong Qin, Yangdong Liu, Yawei Li, Xianglong Liu, Luca Benini, Michele Magno, and Xiaojuan Qi. SliM-LLM: Saliency-driven mixed-precision quantization for large language models. *arXiv:2405.14917*, 2024.
- Itay Hubara, Yury Nahshan, Yair Hanani, Ron Banner, and Daniel Soudry. Accurate post training quantization with small calibration sets. In *International Conference on Machine Learning*, pp. 4466–4475, 2021.
- Hao Kang, Qingru Zhang, Souvik Kundu, Geonhwa Jeong, Zaoxing Liu, Tushar Krishna, and Tuo Zhao. Gear: An efficient kv cache compression recipe for near-lossless generative inference of llm. *arXiv:2403.05527*, 2024.
- Sehoon Kim, Coleman Hooper, Amir Gholami, Zhen Dong, Xiuyu Li, Sheng Shen, Michael W Mahoney, and Kurt Keutzer. SqueezeLLM: Dense-and-sparse quantization. *arXiv:2306.07629*, 2023.

- Changhun Lee, Jungyu Jin, Taesu Kim, Hyungjun Kim, and Eunhyeok Park. OWQ: Outlier-aware weight quantization for efficient fine-tuning and inference of large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 13355–13364, 2024.
- Muyang Li, Yujun Lin, Zhekai Zhang, Tianle Cai, Xiuyu Li, Junxian Guo, Enze Xie, Chenlin Meng, Jun-Yan Zhu, and Song Han. SVDQuant: Absorbing outliers by low-rank components for 4-bit diffusion models. *arXiv:2411.05007*, 2024. URL <https://arxiv.org/abs/2411.05007>.
- Bokai Lin, Zihao Zeng, Zipeng Xiao, Siqi Kou, Tianqi Hou, Xiaofeng Gao, Hao Zhang, and Zhijie Deng. MatryoshkaKV: Adaptive kv compression via trainable orthogonal projection. *arXiv:2410.14731*, 2024a.
- Haokun Lin, Haobo Xu, Yichen Wu, Jingzhi Cui, Yingtao Zhang, Linzhan Mou, Linqi Song, Zhenan Sun, and Ying Wei. Duquant: Distributing outliers via dual transformation makes stronger quantized llms. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024b.
- Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. AWQ: Activation-aware weight quantization for on-device llm compression and acceleration. *Proceedings of Machine Learning and Systems*, 6:87–100, 2024c.
- Yujun Lin, Haotian Tang, Shang Yang, Zhekai Zhang, Guangxuan Xiao, Chuang Gan, and Song Han. QServe: W4a8kv4 quantization and system co-design for efficient llm serving. *arXiv:2405.04532*, 2024d.
- Jing Liu, Ruihao Gong, Xiuying Wei, Zhiwei Dong, Jianfei Cai, and Bohan Zhuang. QLLM: Accurate and efficient low-bitwidth quantization for large language models. *arXiv:2310.08041*, 2023a.
- Tianyang Liu, Canwen Xu, and Julian McAuley. RepoBench: Benchmarking repository-level code auto-completion systems. *arXiv:2306.03091*, 2023b. URL <https://arxiv.org/abs/2306.03091>.
- Zechun Liu, Changsheng Zhao, Igor Fedorov, Bilge Soran, Dhruv Choudhary, Raghuraman Krishnamoorthi, Vikas Chandra, Yuandong Tian, and Tijmen Blankevoort. SpinQuant: Llm quantization with learned rotations. *arXiv:2405.16406*, 2024a.
- Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong, Zhaozhuo Xu, Vladimir Braverman, Beidi Chen, and Xia Hu. KIVI: A tuning-free asymmetric 2bit quantization for kv cache. *arXiv:2402.02750*, 2024b.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. *arXiv:1609.07843*, 2016.
- Meta. Llama 3.2: Revolutionizing edge AI and vision with open, customizable models, 2024a. URL <https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/>.
- Meta. Introducing Meta Llama 3: The most capable openly available LLM to date., 2024b. URL <https://ai.meta.com/blog/meta-llama-3/>.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *EMNLP*, 2018.
- Gunho Park, Baeseong Park, Se Jung Kwon, Byeongwook Kim, Youngjoo Lee, and Dongsoo Lee. nuQmm: Quantized matmul for efficient inference of large-scale generative language models. *arXiv:2206.09557*, 2022.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.

- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. WinoGrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- Charbel Sakr and Brucek Khailany. ESPACE: Dimensionality reduction of activations for model compression. *arXiv:2410.05437*, 2024.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social iqa: Commonsense reasoning about social interactions. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 4463–4473, 2019.
- Utkarsh Saxena, Gobinda Saha, Sakshi Choudhary, and Kaushik Roy. Eigen attention: Attention in low-rank space for KV cache compression. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 15332–15344. Association for Computational Linguistics, 2024. doi: 10.18653/v1/2024.findings-emnlp.899. URL <https://aclanthology.org/2024.findings-emnlp.899>.
- Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang, Peng Gao, Yu Qiao, and Ping Luo. OmniQuant: Omnidirectionally calibrated quantization for large language models. *arXiv:2308.13137*, 2023.
- Sayeh Sharify, Utkarsh Saxena, Zifei Xu, Wanzin Yazar, Ilya Soloveychik, and Xin Wang. Post training quantization of large language models with microscaling formats. In *NeurIPS Efficient Natural Language and Speech Processing Workshop*, pp. 241–258. PMLR, 2024.
- Ying Sheng, Lianmin Zheng, Binhang Yuan, Zhuohan Li, Max Ryabinin, Beidi Chen, Percy Liang, Christopher Ré, Ion Stoica, and Ce Zhang. FlexGen: High-throughput generative inference of large language models with a single gpu. In *International Conference on Machine Learning*, pp. 31094–31116. PMLR, 2023.
- Vijay Thakkar, Pradeep Ramani, Cris Cecka, Aniket Shivam, Honghao Lu, Ethan Yan, Jack Kosian, Mark Hoemmen, Haicheng Wu, Andrew Kerr, Matt Nicely, Duane Merrill, Dustyn Blasig, Fengqi Qiao, Piotr Majcher, Paul Springer, Markus Hohnerbach, Jin Wang, and Manish Gupta. CUTLASS, January 2023. URL <https://github.com/NVIDIA/cutlass>.
- Albert Tseng, Jerry Chee, Qingyao Sun, Volodymyr Kuleshov, and Christopher De Sa. Quip#: Even better llm quantization with hadamard incoherence and lattice codebooks. *arXiv:2402.04396*, 2024. URL <https://arxiv.org/abs/2402.04396>.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-VL: Enhancing vision-language model’s perception of the world at any resolution. *arXiv:2409.12191*, 2024. URL <https://arxiv.org/abs/2409.12191>.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv:1910.03771*, 2020.
- Haocheng Xi, Changhao Li, Jianfei Chen, and Jun Zhu. Training transformers with 4-bit integers. *Advances in Neural Information Processing Systems*, 36:49146–49168, 2023.
- Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*, pp. 38087–38099. PMLR, 2023.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2.5 technical report. *arXiv:2412.15115*, 2024a.

- June Yong Yang, Byeongwook Kim, Jeongin Bae, Beomseok Kwon, Gunho Park, Eunho Yang, Se Jung Kwon, and Dongsoo Lee. No Token Left Behind: Reliable kv cache compression via importance-aware mixed precision quantization. *arXiv:2402.18096*, 2024b.
- Zhewei Yao, Reza Yazdani Aminabadi, Minjia Zhang, Xiaoxia Wu, Conglong Li, and Yuxiong He. ZeroQuant: Efficient and affordable post-training quantization for large-scale transformers. *Advances in Neural Information Processing Systems*, 35:27168–27183, 2022.
- Zhihang Yuan, Lin Niu, Jiawei Liu, Wenyu Liu, Xinggong Wang, Yuzhang Shang, Guangyu Sun, Qiang Wu, Jiayang Wu, and Bingzhe Wu. RPTQ: Reorder-based post-training quantization for large language models. *arXiv:2304.01089*, 2023a.
- Zhihang Yuan, Yuzhang Shang, Yue Song, Qiang Wu, Yan Yan, and Guangyu Sun. ASVD: Activation-aware singular value decomposition for compressing large language models. *arXiv:2312.05821*, 2023b.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhao Chen. MMMU: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. *arXiv:2311.16502*, 2024. URL <https://arxiv.org/abs/2311.16502>.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a machine really finish your sentence? *arXiv:1905.07830*, 2019.
- Chao Zeng, Songwei Liu, Yusheng Xie, Hong Liu, Xiaojian Wang, Miao Wei, Shu Yang, Fangmin Chen, and Xing Mei. ABQ-LLM: Arbitrary-bit quantized inference acceleration for large language models. *arXiv:2408.08554*, 2024.
- Yilong Zhao, Chien-Yu Lin, Kan Zhu, Zihao Ye, Lequn Chen, Size Zheng, Luis Ceze, Arvind Krishnamurthy, Tianqi Chen, and Baris Kasikci. Atom: Low-bit quantization for efficient and accurate llm serving. *Proceedings of Machine Learning and Systems*, 6:196–209, 2024.

## A RELATED WORKS

### A.1 QUANTIZATION OF LLMs

Quantization reduces model size and accelerates inference by lowering neural network bit precision (Choi et al., 2018; Hubara et al., 2021; Yao et al., 2022; Park et al., 2022; Gholami et al., 2022; Xi et al., 2023). It is broadly categorized into two categories: *uniform precision quantization* (UPQ) and *mixed precision quantization* (MPQ). **Uniform precision quantization (UPQ)** applies the same bit-width across all layers, simplifying implementation but neglecting layer-specific sensitivity to quantization. **Weight-only UPQ** methods reduce storage by compressing weights, using techniques like Hessian-guided rounding (GPTQ, Frantar et al. 2022), adaptive rounding (QuIP, Chee et al. 2024), channel-wise scaling (AWQ, Lin et al. 2024c), and multi-codebook quantization (AQLM, Egiazarian et al. 2024). However, these methods struggle with batch processing due to significant activation memory overhead. **Weight-activation UPQ** compresses both weights and activations to address this. Methods such as SmoothQuant (Xiao et al., 2023) and OmniQuant (Shao et al., 2023) scale activations and weights to handle outliers, while RPTQ (Yuan et al., 2023a), QLLM (Liu et al., 2023a), and QServe (Lin et al., 2024d) employ channel-level strategies like clustering and reordering. Rotation-based methods such as QuaRot (Ashkboos et al., 2024c), SpinQuant (Liu et al., 2024a) and DuQuant Lin et al. (2024b) further enhance robustness in low-precision scenarios. **KV cache UPQ** reduces memory for large batches or long contexts. FlexGen (Sheng et al., 2023) employs 4-bit quantization and memory offloading, while KIVI (Liu et al., 2024b) uses asymmetric 2-bit quantization for compression, enabling efficient inference.

**Mixed precision quantization (MPQ)** optimizes bit-widths by adapting to the sensitivity of weights and activations, achieving better accuracy than UPQ at similar compression rates. *Our proposed method, ResQ, follows the MPQ approach.* **Weight-only MPQ** has advanced efficiency for memory-bound applications with minimal activation demands. Methods like OWQ (Lee et al., 2024) and SpQR (Dettmers et al., 2023) mitigate activation outliers’ impact by retaining critical features in full precision, while SqueezeLLM (Kim et al., 2023) employs Dense-and-Sparse decomposition to efficiently store sensitive weights. **Weight-activation MPQ** enhances efficiency by addressing activation outliers (e.g. Guan et al. 2024; Zeng et al. 2024). Methods like LLM.int8() (Dettmers et al., 2022) and QUIK(Ashkboos et al., 2024b) preserve critical activations with mixed or low-precision decompositions, while Atom (Zhao et al., 2024) and Slim-LLM (Huang et al., 2024) optimize quantization through channel reordering and salience-driven bit allocation. **KV cache MPQ** reduces memory usage while preserving precision for critical tokens using techniques like non-uniform quantization, importance-aware precision, and salient token compression (Hooper et al., 2024; Yang et al., 2024b; Dong et al., 2024; He et al., 2024). Alternatively, GEAR quantizes all tokens’ KV cache and maintains low-rank quantization error (Kang et al., 2024).

## A.2 LOW-RANK DECOMPOSITION

Low-rank decomposition techniques have been widely used in model compression, reducing dimensionality while maintaining performance. For instance, SliceGPT Ashkboos et al. (2024a) projects weight matrices onto principal components for sparsification, while ESPACE Sakr & Khailany (2024) reduces activation dimensionality via pre-calibrated projections, achieving inference-time efficiency. Similarly, ASVD Yuan et al. (2023b) introduces an activation-aware decomposition method that incorporates activation distributions into weight decomposition. Additionally, low-rank decomposition can be applied to reduce KV cache size. For example, Eigen Attention Saxena et al. (2024) and ASVD Yuan et al. (2023b) employ low-rank approximations to reduce memory usage in KV caches during attention operations. PALU Chang et al. (2024) introduces learnable projections to adaptively compress KV caches based on the compression budget. Finally, Matryoshka KV Cache refines this with hierarchical orthogonal projections and knowledge distillation Lin et al. (2024a).

## B PROOF OF THEOREM 2.1

We begin the proof by introducing the following lemma.

**Lemma B.1.** *For any tensor  $\mathbf{R}$  quantized following the quantization described in equation ??, assuming the values of  $\mathbf{R}$  follows a normal distribution, we have*

$$\mathbb{E}\|\mathbf{R} - Q(\mathbf{R})\|_F \leq \frac{\sqrt{\pi \log [\text{size}(\mathbf{R})]}}{2^{n-1} - 1} \mathbb{E}\|\mathbf{R}\|_F \quad (5)$$

where  $\text{size}(\mathbf{R})$  denotes the number of elements in  $\mathbf{R}$ .

Proof of lemma B.1 can be found in Li et al. (2024). From this lemma we obtain that the quantization error  $\|\mathbf{R} - Q(\mathbf{R})\|_F$  is bounded by the magnitude of the tensor quantized  $\|\mathbf{R}\|_F$ . Now for our use case of mixed precision quantization where the low-precision component is quantized to  $L$  bits and high precision component is quantized to  $H$  bits, we write the quantization error again below,

$$\begin{aligned} \mathbb{E}\|\mathbf{X} - \mathbf{X}_q\|_F &= \mathbb{E}\|\mathbf{X}\mathbf{U}_l - Q_L(\mathbf{X}\mathbf{U}_l)\|_F \\ &+ \mathbb{E}\|\mathbf{X}\mathbf{U}_h - Q_H(\mathbf{X}\mathbf{U}_h)\|_F. \end{aligned} \quad (6)$$

The random rotation matrices  $\mathbf{R}$  ensure that  $\mathbf{X}\mathbf{U}_l$  and  $\mathbf{X}\mathbf{U}_h$  are normally distributed by Lemma ???. Applying Lemma B.1 to the quantization error in equation 6, we get,

$$\begin{aligned}
\|\mathbf{X} - \mathbf{X}_q\|_F &\leq \frac{\sqrt{\log(\text{size}(\mathbf{X}\mathbf{U}_l))\pi}}{2^{L-1} - 1} \mathbb{E}\|\mathbf{X}\mathbf{U}_l\|_F \\
&+ \frac{\sqrt{\log(\text{size}(\mathbf{X}\mathbf{U}_h))\pi}}{2^{H-1} - 1} \mathbb{E}\|\mathbf{X}\mathbf{U}_h\|_F \\
&= \frac{\sqrt{\log(\text{size}(\mathbf{X}\mathbf{P}_l))\pi}}{2^{L-1} - 1} \mathbb{E}\|\mathbf{X}\mathbf{P}_l\|_F \\
&+ \frac{\sqrt{\log(\text{size}(\mathbf{X}\mathbf{P}_h))\pi}}{2^{H-1} - 1} \mathbb{E}\|\mathbf{X}\mathbf{P}_h\|_F \\
&= \frac{\sqrt{\log(\text{size}(\mathbf{X}\mathbf{P}_l))\pi}}{2^{L-1} - 1} \mathbb{E}\|\text{tr}(\mathbf{X}\mathbf{P}_l\mathbf{P}_l^\top \mathbf{X}^\top)\|_F \\
&+ \frac{\sqrt{\log(\text{size}(\mathbf{X}\mathbf{P}_h))\pi}}{2^{H-1} - 1} \mathbb{E}\|\text{tr}(\mathbf{X}\mathbf{P}_h\mathbf{P}_h^\top \mathbf{X}^\top)\|_F
\end{aligned} \tag{7}$$

We know  $\text{size}(\mathbf{X}\mathbf{P}_l) = d - r$  and  $\text{size}(\mathbf{X}\mathbf{P}_h) = r$  since  $r$  components are in high precision. With  $\mathbf{P}_l\mathbf{P}_l^\top + \mathbf{P}_h\mathbf{P}_h^\top = \mathbf{I}$ , we have

$$\begin{aligned}
\|\mathbf{X} - \mathbf{X}_q\|_F &\leq \frac{\sqrt{\log(d-r)\pi}}{2^{L-1} - 1} (\mathbb{E}\|\mathbf{X}\|_F - \mathbb{E}\|\mathbf{X}\mathbf{P}_h\|_F) \\
&+ \frac{\sqrt{\log(r)\pi}}{2^{H-1} - 1} \mathbb{E}\|\mathbf{X}\mathbf{P}_h\|_F \\
&= \frac{\sqrt{\log(d-r)\pi}}{2^{L-1} - 1} \mathbb{E}\|\mathbf{X}\|_F \\
&- \left( \frac{\sqrt{\log(d-r)\pi}}{2^{L-1} - 1} - \frac{\sqrt{\log(r)\pi}}{2^{H-1} - 1} \right) \mathbb{E}\|\mathbf{X}\mathbf{P}_h\|_F
\end{aligned} \tag{8}$$

Since  $\frac{\sqrt{\log(d-r)\pi}}{2^{L-1} - 1} - \frac{\sqrt{\log(r)\pi}}{2^{H-1} - 1} > 0$  the quantization error is reduced by maximizing  $\|\mathbf{X}\mathbf{P}_h\|_F$

## C ADDITIONAL IMPLEMENTATION DETAILS

We implement ResQ using the HuggingFace Transformers library Wolf et al. (2020) with PyTorch Paszke et al. (2019). We share a single  $\mathbf{U}_A$  across all layers, while  $\mathbf{U}_B$ ,  $\mathbf{U}_C$  and  $\mathbf{U}_D$  are generated per layer. Following SpinQuant Liu et al. (2024a), we use per-token asymmetric quantization for activations, per-channel symmetric quantization for weights, and per-head asymmetric quantization for the KV cache. We fuse the projection matrices  $\mathbf{U}_A$ ,  $\mathbf{U}_B$ ,  $\mathbf{U}_D$  into weights and apply GPTQ Frantar et al. (2022) for weight quantization. To efficiently implement on-the-fly projections,  $\mathbf{U}_D$  is a Hadamard matrix and  $\mathbf{U}_C$  and its activations are quantized to 8-bit.

In this work, obtaining the projection matrices and quantization of weights for **all the models** is performed on a single NVIDIA A100 80GB GPUs. Time taken by ResQ compared with other approaches is shown in Table ??. Evaluation on various benchmarks for all the models is also done on a single NVIDIA A100 GPU with the sole exception of Meta-Llama-3-70b which requires 4 GPUs for evaluation. We use `lm.evaluation.harness` version 0.4.5 Gao et al. (2024) and LongBench Bai et al. (2024) for all the evaluation tasks. For Arc-c/e, Hellaswag, OpenBook QA, PIQA tasks we report `acc_norm` while for BoolQ, SIQA and Winogrande we report `acc`.

For calibration data, we use 512 randomly chosen samples for Wikitext to obtain the projection matrices. While for GPTQ we use 128 randomly chosen samples from Wikitext following the original work Frantar et al. (2022).

The KV cache, as well as the weights and activations of all Linear layers (except `mlp.down_proj`), are quantized to 4-bit precision, with  $\frac{1}{8}$  of channels retained in 8-bit precision. While, the weights and activations within `down_proj` are uniformly quantized to 4-bit precision. Following Ashkboos et al. (2024c) and Liu et al. (2024a), we keep query vector in 16-bit.

## D COMPLETE RESULTS OF MAIN RESULT TABLES

Detailed results of Table 1 in the main paper, including more models and task-by-task performance, are shown in Tables 4 (Llama families) and 5 (Qwen2.5 family). As expected, ResQ achieves superior performance to baselines across the series of common sense reasoning and MMLU tasks.

Table 4: Comparison of perplexity on Wikitext, accuracy on eight 0-shot common sense reasoning tasks including ARC-challenge, ARC-easy, BoolQ, HellaSwag, Openbook QA, PIQA, SIQA, and WinoGrande, and 0-shot massive multitask language understanding tasks across four subjects: STEM, Humanities, Social Sciences, and MMLU-other, for the Llama 2, Llama 3 and Llama 3.2 families when quantized to  $W/A/KV = 4/4/4$  bits. Results of all techniques were obtained using their official codebase. \*: Mixed precision with  $1/8$  channels in 8-bit for  $W/A/KV = 4.5$ -bit. All techniques except RTN use GPTQ Frantar et al. (2022). ( $\downarrow$ ): lower is better, ( $\uparrow$ ): higher is better.

Llama 2 family																
Model	Method	0-shot common sense reasoning tasks										0-shot MMLU tasks				
		Perplexity (↓)	Wiki (↑)	ARC-c (↑)	ARC-e (↑)	BoolQ (↑)	HellaS (↑)	OBQA (↑)	PIQA (↑)	SIQA (↑)	WinoG (↑)	Avg. (↑)	humanities (↑)	Other (↑)	SocialS (↑)	STEM (↑)
Llama-2-7b-hf	16-bit	5.5	46.3	74.6	77.8	75.9	44.2	79.2	46.1	69.1	64.1	38.9	43.9	46.0	33.4	41.1
	RTN	1766.2	26.3	27.8	54.8	29.4	23.8	31.0	35.0	48.7	37.4	24.5	24.7	22.9	22.2	23.6
	GPTQ	9600.0	24.8	31.4	55.4	30.6	25.6	55.8	34.2	53.3	38.9	24.7	24.5	22.7	23.2	23.8
	SmoothQuant+	15.4	29.3	47.1	56.8	48.6	31.8	65.5	37.2	52.4	46.1	25.0	24.5	24.1	23.4	24.2
	QUIK*	7.5	39.8	63.7	68.9	68.3	37.8	72.9	42.1	62.4	57.0	26.9	29.6	28.8	25.8	27.8
	QuaRot	6.1	41.5	71.4	73.2	73.2	40.6	76.9	43.6	65.6	60.7	31.2	35.1	34.6	28.2	32.3
	SpinQuant	6.0	43.6	71.3	73.8	73.2	40.4	76.0	44.1	65.4	61.0	33.9	38.5	37.5	29.5	34.8
ResQ*	5.8	44.0	72.6	75.3	74.0	41.0	77.9	43.9	66.9	62.0	35.9	40.9	42.2	32.2	37.7	
Llama-2-13b-hf	16-bit	4.9	49.1	77.4	80.5	79.4	45.2	80.7	47.2	72.1	66.5	47.9	59.3	61.0	42.4	52.7
	RTN	3543.9	22.8	29.8	40.2	26.6	27.8	31.4	35.6	50.6	33.5	23.7	25.0	23.1	22.6	23.6
	GPTQ	3120.0	23.6	31.1	38.7	27.2	26.8	53.6	35.8	49.8	33.8	25.0	25.4	23.7	25.1	24.8
	SmoothQuant+	11.2	34.5	55.6	62.9	62.5	32.4	70.1	38.7	55.6	51.0	25.7	26.1	27.3	26.6	26.6
	QUIK*	6.8	43.7	68.0	71.3	73.3	40.0	75.7	45.1	64.6	60.2	34.7	40.6	39.8	31.8	36.7
	QuaRot	5.4	46.9	74.9	76.6	75.8	42.6	79.1	45.5	69.0	63.8	43.8	53.6	54.0	39.4	47.7
	SpinQuant	5.2	49.0	76.3	78.2	77.1	42.8	79.3	46.3	69.5	64.8	43.5	53.1	55.4	39.1	47.8
ResQ*	5.1	49.1	76.1	79.7	77.9	43.6	79.1	46.6	69.9	65.2	45.3	56.0	58.0	41.0	50.1	
Llama 3 family																
Model	Method	0-shot common sense reasoning tasks										0-shot MMLU tasks				
		Perplexity (↓)	Wiki (↑)	ARC-c (↑)	ARC-e (↑)	BoolQ (↑)	HellaS (↑)	OBQA (↑)	PIQA (↑)	SIQA (↑)	WinoG (↑)	Avg. (↑)	humanities (↑)	Other (↑)	SocialS (↑)	STEM (↑)
Meta-Llama-3-8B	16-bit	6.1	53.2	77.1	81.1	79.2	44.8	80.9	47.0	73.4	67.1	55.0	70.6	73.2	53.7	63.1
	RTN	218.9	25.3	34.9	44.2	38.3	27.8	56.5	36.8	50.8	39.3	24.7	25.1	23.3	21.4	23.6
	GPTQ	166.3	24.7	37.7	44.3	36.8	27.0	57.6	36.4	53.8	39.8	24.7	23.9	22.8	21.8	23.3
	SmoothQuant+	78.2	27.5	42.0	50.7	44.9	28.8	59.0	35.9	50.9	42.5	25.4	25.5	24.5	23.4	24.7
	QUIK*	14.2	33.6	56.4	60.5	61.5	33.2	68.7	39.9	59.0	51.6	30.0	34.0	34.8	32.1	32.7
	QuaRot	7.8	45.1	70.4	73.8	74.7	42.6	76.6	45.1	68.5	62.1	47.8	59.1	61.4	44.3	53.2
	SpinQuant	7.4	48.0	75.4	75.8	75.4	43.8	77.5	45.0	69.2	63.8	49.8	63.3	65.0	46.8	56.2
ResQ*	7.1	49.2	75.0	72.5	76.5	43.0	78.3	45.8	71.0	63.9	50.6	64.4	65.8	48.1	57.2	
Meta-Llama-3-70B	16-bit	2.9	64.2	85.9	85.3	84.9	48.6	84.4	50.8	86.6	73.1	67.6	81.5	86.8	68.4	76.1
	RTN	452.7	32.6	50.3	54.2	41.3	31.6	64.8	35.9	53.2	45.5	24.5	25.8	22.3	22.1	25.2
	GPTQ	11655.0	25.9	26.0	37.9	26.2	28.6	50.4	34.3	49.9	34.9	27.1	24.3	24.0	26.5	25.5
	SmoothQuant+	8.0	44.5	68.9	60.7	75.0	36.4	76.1	43.2	60.4	58.2	46.6	56.4	58.0	43.6	51.1
	QUIK*	5.7	53.7	74.5	81.6	81.1	46.6	81.0	46.8	75.2	67.6	55.7	72.5	75.8	57.3	65.3
	QuaRot	6.2	52.0	77.3	81.7	75.6	43.8	78.8	43.4	72.8	65.7	50.7	67.0	68.1	51.9	59.4
	SpinQuant	4.1	61.4	84.3	83.9	83.5	46.0	83.1	48.6	78.3	71.1	64.9	79.9	84.9	66.1	74.0
ResQ*	4.1	61.4	84.3	83.9	83.5	46.0	83.1	48.6	78.3	71.1	64.9	79.9	84.9	66.1	74.0	
Llama 3.2 family																
Model	Method	0-shot common sense reasoning tasks										0-shot MMLU tasks				
		Perplexity (↓)	Wiki (↑)	ARC-c (↑)	ARC-e (↑)	BoolQ (↑)	HellaS (↑)	OBQA (↑)	PIQA (↑)	SIQA (↑)	WinoG (↑)	Avg. (↑)	humanities (↑)	Other (↑)	SocialS (↑)	STEM (↑)
Llama-3.2-1B	16-bit	9.8	36.5	60.6	63.4	63.6	37.4	74.5	42.8	60.1	54.9	34.8	41.1	39.9	32.0	36.9
	RTN	329.1	22.4	29.9	53.4	31.4	29.4	54.8	34.9	48.5	38.1	24.8	25.2	22.4	22.7	23.8
	GPTQ	108.9	24.7	32.7	52.3	30.7	23.6	54.3	34.4	51.1	38.0	24.7	25.1	25.5	24.5	24.9
	SmoothQuant+	228.9	23.3	30.1	52.9	31.3	26.6	54.2	34.5	51.2	38.0	23.9	24.1	25.0	23.5	24.1
	QUIK*	21.8	27.4	46.0	55.0	46.0	26.4	62.4	38.6	52.6	44.3	25.6	25.6	24.6	24.5	25.1
	QuaRot	14.3	30.0	51.4	59.1	54.0	34.2	66.7	39.6	57.1	49.0	25.4	26.9	25.4	24.4	25.5
	SpinQuant	13.6	32.3	51.8	59.3	55.4	30.4	67.7	38.6	54.7	48.8	25.4	27.6	24.2	25.3	25.6
ResQ*	12.4	34.0	54.2	57.0	57.3	31.2	69.4	41.0	56.8	50.1	28.3	30.5	31.3	27.6	29.4	
Llama-3.2-3B	16-bit	7.8	46.2	71.7	73.1	73.7	43.4	77.4	47.2	69.1	62.7	48.9	62.9	62.3	45.2	54.8
	RTN	268.8	23.5	35.4	46.2	33.6	28.2	36.3	33.6	50.6	35.7	25.1	25.6	27.0	24.9	25.7
	GPTQ	178.3	27.0	27.0	48.8	44.4	27.8	59.1	37.1	51.5	40.3	24.9	24.5	25.7	24.0	24.8
	SmoothQuant+	96.1	25.3	33.1	47.8	37.7	25.2	56.2	35.8	50.9	39.0	25.4	26.6	26.4	25.3	25.9
	QUIK*	15.8	32.9	50.1	52.6	59.1	33.2	68.7	40.3	53.0	48.8	29.0	33.2	31.9	30.3	31.1
	QuaRot	10.1	38.6	59.0	65.9	66.5	35.8	74.4	43.1	65.2	56.1	38.5	47.3	46.7	35.3	42.0
	SpinQuant	9.2	38.9	64.8	68.0	69.1	39.4	74.9	45.1	62.9	57.9	37.0	49.4	50.5	39.9	44.2
ResQ*	8.8	43.1	65.6	68.8	70.5	38.4	75.1	45.6	64.8	59.0	44.7	57.0	56.5	41.0	49.8	

## E COMPLETE RESULTS OF THE MMMU BENCHMARK

This section presents task-by-task results for the MMMU benchmark across six subjects—Art & Design, Business, Science, Health & Medicine, Humanities & Social Science, and Tech & Engineering—for the Qwen2 VL family when quantized to  $W/A/KV = 4/4/4$  bits and  $W/A/KV = 4/8/4$  bits of precision. On average, ResQ consistently outperforms all baselines across different models. Notably, the advantage of ResQ becomes more pronounced with larger models. For instance, for Qwen2-VL-7B-Instruct at  $W/A/KV = 4/8/4$  bits of precision, ResQ achieves an average accuracy score of 48.8, significantly outperforming the next-best method, QUIK, which scores 26.4, representing an  $\sim 85\%$  relative improvement.

Table 5: Comparison of perplexity score on Wikitext, accuracy on eight 0-shot common sense reasoning tasks including ARC-challenge, ARC-easy, BoolQ, HellaSwag, Openbook QA, PIQA, SIQA, and WinoGrande, and 0-shot massive multitask language understanding tasks across four subjects: STEM, Humanities, Social Sciences, and MMLU-other, for the Qwen2.5 family when quantized to W/A/KV = 4/4/4 bits. Results of all techniques were obtained using their official code-base. \*: Mixed precision with 1/8 channels in 8-bit for W/A/KV = 4.5-bit. All techniques except RTN use GPTQ Frantar et al. (2022). (↓): lower is better, (↑): higher is better.

		Qwen2.5 family														
Model	Method	Perplexity	0-shot common sense reasoning tasks								0-shot MMLU tasks					
		Wiki (↓)	ARC-c (↑)	ARC-e (↑)	BoolQ (↑)	HellaS (↑)	OBQA (↑)	PIQA (↑)	SIQA (↑)	WinoG (↑)	Avg. (↑)	humanities (↑)	Other (↑)	SocialS (↑)	STEM (↑)	Avg. (↑)
Qwen2.5-0.5B	16-bit	13.1	31.9	58.4	62.1	52.1	35.0	69.7	44.3	57.1	51.3	42.2	53.2	55.5	41.5	48.1
	RTN	23204.3	26.2	27.0	39.3	26.0	24.0	50.7	34.5	51.5	34.9	24.8	24.0	22.8	24.3	25.9
	GPTQ	16302.3	23.7	26.9	39.0	26.5	26.4	50.2	33.4	49.6	34.5	24.1	24.8	23.5	23.0	23.9
	SmoothQuant+	10053.9	25.9	26.3	39.9	27.2	25.4	47.1	35.9	49.6	34.7	24.5	24.7	21.5	22.1	23.2
	QUIK*	38.6	24.5	38.6	48.0	36.9	28.4	58.1	36.4	51.9	40.4	26.3	25.9	23.6	24.2	25.0
	QuaRot	219.9	25.4	36.6	45.0	28.9	28.6	54.1	32.9	51.7	37.9	24.4	24.0	23.0	23.5	23.7
ResQ*	<b>29.6</b>	<b>27.1</b>	<b>44.2</b>	<b>53.2</b>	<b>38.8</b>	<b>28.0</b>	<b>61.9</b>	<b>34.4</b>	<b>51.3</b>	<b>42.4</b>	<b>26.1</b>	<b>27.5</b>	<b>25.3</b>	<b>26.0</b>	<b>26.2</b>	
Qwen2.5-1.5B	16-bit	9.3	45.1	72.1	72.9	67.7	40.2	76.3	48.8	63.7	60.8	53.5	65.5	70.6	52.8	60.6
	RTN	14518.9	23.1	27.2	43.9	26.8	25.6	51.5	33.4	52.5	35.5	23.8	24.5	23.8	22.7	25.7
	GPTQ	25769.7	23.9	26.9	43.9	26.1	27.6	49.7	32.1	51.5	35.2	24.6	24.7	23.7	23.8	24.2
	SmoothQuant+	31655.9	25.0	26.2	39.9	26.0	26.0	50.8	32.1	49.0	34.4	25.5	24.4	22.7	22.4	23.8
	QUIK*	6613.5	21.8	31.9	40.9	27.9	27.4	52.8	35.2	48.6	35.8	24.6	24.0	21.9	21.7	23.1
	QuaRot	6599.9	23.6	37.3	46.2	28.6	27.0	56.3	35.2	52.4	38.3	24.5	24.3	23.0	22.4	23.5
ResQ*	<b>12.5</b>	<b>38.7</b>	<b>64.1</b>	<b>65.7</b>	<b>61.4</b>	<b>37.8</b>	<b>71.6</b>	<b>42.7</b>	<b>60.1</b>	<b>55.3</b>	<b>43.2</b>	<b>54.4</b>	<b>54.9</b>	<b>41.5</b>	<b>48.5</b>	
Qwen2.5-3B	16-bit	8.0	47.4	73.0	77.5	73.6	42.0	78.7	49.9	68.4	63.8	56.6	71.0	76.3	60.6	66.1
	RTN	39033.0	25.6	25.8	41.7	26.3	27.4	49.5	33.1	51.4	35.1	24.5	24.4	22.8	21.9	23.4
	GPTQ	9977.8	26.0	26.7	41.5	26.7	28.2	51.5	31.9	48.3	35.1	24.3	23.8	22.8	21.8	23.2
	SmoothQuant+	73306.7	25.4	24.5	41.0	26.4	29.8	48.4	32.4	50.4	34.8	25.6	24.7	23.1	22.4	23.9
	QUIK*	15.5	36.1	55.4	61.4	57.2	36.2	67.1	40.8	55.3	51.2	36.4	42.8	42.4	36.1	39.4
	QuaRot	68.8	32.4	53.1	51.6	49.2	33.4	66.7	39.3	56.4	47.7	28.1	32.0	28.9	26.6	28.9
ResQ*	<b>9.0</b>	<b>45.3</b>	<b>70.5</b>	<b>72.7</b>	<b>70.2</b>	<b>42.4</b>	<b>76.8</b>	<b>46.7</b>	<b>64.4</b>	<b>61.1</b>	<b>53.1</b>	<b>66.5</b>	<b>70.5</b>	<b>54.8</b>	<b>61.2</b>	
Qwen2.5-7B	16-bit	6.8	51.2	77.6	84.7	78.9	47.2	80.0	54.8	73.2	68.4	62.6	76.7	82.6	70.1	73.0
	RTN	24382.1	24.5	26.3	37.8	26.0	29.0	51.0	34.1	50.1	34.9	24.9	24.5	23.4	24.9	24.4
	GPTQ	13593.7	25.2	25.6	37.8	26.3	28.2	52.4	34.4	48.9	34.8	24.4	24.3	22.8	22.6	23.5
	SmoothQuant+	19088.7	26.3	25.2	39.8	26.4	27.6	52.7	33.5	52.0	35.4	25.1	25.4	22.6	24.1	24.3
	QUIK*	260.3	29.5	42.4	51.7	36.3	28.2	59.6	34.5	49.6	41.5	24.3	26.9	23.1	23.8	24.6
	QuaRot	4035.9	25.9	41.0	39.1	29.1	27.6	57.9	35.7	50.6	38.4	24.8	24.4	24.4	22.7	24.1
ResQ*	<b>8.2</b>	<b>49.0</b>	<b>74.7</b>	<b>81.4</b>	<b>75.7</b>	<b>45.0</b>	<b>78.9</b>	<b>49.4</b>	<b>68.2</b>	<b>65.3</b>	<b>57.8</b>	<b>74.4</b>	<b>79.3</b>	<b>64.5</b>	<b>69.0</b>	
Qwen2.5-14B	16-bit	5.3	58.8	79.4	85.4	82.9	45.4	81.9	55.3	75.8	70.6	69.9	81.9	86.2	76.5	78.6
	RTN	2715	21.6	32.7	51.5	29.6	25.8	52.6	33.2	51.7	37.3	25.3	25.2	26.0	25.3	24.9
	GPTQ	5100.3	23.8	29.1	47.7	30.1	27.6	51.3	34.6	51.2	36.9	25.1	24.7	25.1	24.3	24.8
	SmoothQuant+	1375.7	27.0	26.3	38.0	26.8	29.2	51.6	32.4	49.3	35.1	25.9	24.5	22.2	22.2	23.7
	QUIK*	10.5	45.0	67.1	64.7	68.9	37.6	74.8	43.9	59.3	57.6	48.9	61.1	64.7	51.5	56.6
	QuaRot	6.8	54.8	79.6	79.9	78.7	44.0	79.5	49.9	70.7	67.1	60.9	75.1	80.2	67.3	70.9
ResQ*	<b>6.2</b>	<b>57.6</b>	<b>82.1</b>	<b>84.9</b>	<b>81.1</b>	<b>44.8</b>	<b>80.5</b>	<b>51.7</b>	<b>70.6</b>	<b>69.2</b>	<b>65.2</b>	<b>78.4</b>	<b>83.4</b>	<b>71.5</b>	<b>74.6</b>	
Qwen2.5-32B	16-bit	5.0	55.7	78.0	87.4	84.1	44.4	82.3	56.4	75.2	70.4	73.1	83.6	89.6	81.2	81.9
	RTN	1847.4	24.3	35.3	51.4	31.9	27.0	52.8	34.1	51.4	38.5	24.5	25.1	25.3	24.3	24.8
	GPTQ	3891.1	25.4	35.4	48.5	31.8	27.0	53.8	35.8	50.5	38.5	25.9	24.8	23.6	24.0	24.6
	SmoothQuant+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	QUIK*	9.6	41.0	64.6	74.9	72.0	39.6	75.8	44.5	60.2	59.1	54.7	66.8	71.3	58.8	62.9
	QuaRot	6.1	54.5	76.1	85.1	81.5	44.2	80.1	51.3	70.4	67.9	68.5	80.0	86.0	76.0	77.6
ResQ*	<b>5.6</b>	<b>55.1</b>	<b>78.4</b>	<b>86.0</b>	<b>82.5</b>	<b>45.4</b>	<b>81.1</b>	<b>53.9</b>	<b>74.0</b>	<b>69.5</b>	<b>70.3</b>	<b>82.3</b>	<b>87.9</b>	<b>78.9</b>	<b>79.8</b>	
Qwen2.5-72B	16-bit	3.9	62.6	83.2	89.2	86.0	46.6	83.6	58.4	77.7	73.4	77.2	86.9	90.6	82.4	84.3
	RTN	45412.7	25.9	26.3	38.0	25.9	25.2	50.0	34.2	48.7	34.3	25.5	24.2	23.0	23.2	24.0
	GPTQ	37967.2	25.4	25.8	38.1	25.6	26.6	51.2	34.2	49.4	34.5	25.1	24.0	21.9	22.2	23.3
	SmoothQuant+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	QUIK*	8.3	45.1	68.1	77.2	77.2	39.0	77.4	45.6	65.6	61.9	60.2	74.3	77.5	65.3	69.3
	QuaRot	4.9	55.8	81.1	87.5	84.0	45.2	81.7	52.5	74.5	70.3	71.4	84.2	87.7	77.1	80.1
ResQ*	<b>4.6</b>	<b>58.4</b>	<b>80.9</b>	<b>88.4</b>	<b>84.9</b>	<b>48.2</b>	<b>82.6</b>	<b>55.5</b>	<b>77.0</b>	<b>72.0</b>	<b>72.8</b>	<b>84.6</b>	<b>89.0</b>	<b>79.5</b>	<b>81.5</b>	

Table 6: Accuracy (higher is better) on 0-shot massive multi-discipline multimodal understanding and reasoning tasks across six subjects: Art & Design, Business, Science, Health & Medicine, Humanities & Social Science, and Tech & Engineering for the Qwen2 VL Instruct family. Results of all techniques were obtained using their official codebase. Our work ResQ and QUIK Ashkboos et al. (2024b) keep  $1/8$  of channels in 8-bit for average of 4.5-bit. All techniques except RTN use GPTQ Frantar et al. (2022).

Qwen2-VL-2B-Instruct								
W/A/KV (bit)	Method	0-shot MMMU tasks						
		Art-Design	Business	Science	Health	Humanities	Tech	Avg.
16/16/16	Baseline	56.7	36.0	37.3	50.8	26.0	31.0	39.6
	RTN	28.3	18.7	26.0	26.7	21.3	<b>29.1</b>	25.0
4/4/4	GPTQ	28.3	<b>27.3</b>	27.0	29.0	<b>26.7</b>	27.6	27.7
	QuaRot	24.2	23.3	20.7	26.7	26.0	22.9	24.0
	QUIK	25.8	26.0	26.7	29.2	26.0	24.3	26.3
4.5/4.5/4.5	ResQ	<b>38.3</b>	21.3	<b>28.7</b>	<b>45.0</b>	21.3	23.3	<b>29.7</b>
4/8/4	RTN	27.5	21.3	27.3	24.2	21.3	27.6	24.9
	GPTQ	24.2	23.3	24.0	18.3	21.3	29.5	23.4
	QuaRot	20.0	24.7	30.0	26.7	26.0	<b>31.4</b>	26.5
	QUIK	33.3	28.7	32.0	32.5	26.0	18.1	28.4
4.5/8/4.5	ResQ	<b>37.5</b>	<b>32.0</b>	<b>32.7</b>	<b>47.5</b>	<b>26.7</b>	27.6	<b>34.0</b>

Qwen2-VL-7B-Instruct								
W/A/KV (bit)	Method	0-shot MMMU tasks						
		Art-Design	Business	Science	Health	Humanities	Tech	Avg.
16/16/16	Baseline	68.3	41.3	54.7	68.3	38.7	38.1	51.6
	RTN	24.2	28.0	29.3	22.5	29.3	27.1	26.7
4/4/4	GPTQ	21.7	26.0	25.3	28.3	24.7	23.3	24.9
	QuaRot	21.7	21.3	28.7	25.0	20.7	29.5	24.5
	QUIK	30.8	30.0	32.0	26.7	28.0	26.2	28.9
4.5/4.5/4.5	ResQ	<b>65.0</b>	<b>39.3</b>	<b>45.3</b>	<b>61.7</b>	<b>34.0</b>	<b>36.7</b>	<b>47.0</b>
4/8/4	RTN	23.3	28.7	27.3	25.0	22.7	24.3	25.2
	GPTQ	20.8	23.3	30.0	19.2	24.0	28.6	24.3
	QuaRot	20.8	26.0	30.0	19.2	24.7	26.2	24.5
	QUIK	25.0	23.3	31.3	26.7	25.3	26.7	26.4
4.5/8/4.5	ResQ	<b>67.5</b>	<b>39.3</b>	<b>51.3</b>	<b>64.2</b>	<b>36.7</b>	<b>33.8</b>	<b>48.8</b>