

MINSOO KIM

222, Wangsimni-ro, Seongdong-gu, Seoul, Republic of Korea, 04763

☎ +82-10-8203-6871 🏠 marsjacobs.github.io ✉ minsoo2333@hanyang.ac.kr

RESEARCH INTERESTS

General Efficiency for LLMs - quantization, distillation, parameter-efficient fine-tuning, kv-cache compression

EDUCATION

Hanyang University, Seoul, South Korea

Mar. 2021 - Present

Artificial Intelligence Hardware & Algorithm lab

Ph.D. Student in Electronic Engineering

Advisor: Professor Jungwook Choi

Hanyang University, Seoul, South Korea

Feb. 2021

B.S in Electronic Engineering

Thesis: Improving training method for very low bit weight quantization of Light Deep Learning Model

Advisor: Professor Jungwook Choi

WORK EXPERIENCE

Qualcomm AI Research, Seoul, South Korea, PhD research Intern

Mar. 2024 - Sep. 2024

Hanyang University, Seoul, South Korea, Student researcher

Feb. 2021 - Present

PUBLICATIONS

- **[ACL 2024]** Minsoo Kim, Sihwa Lee, Wonyong Sung and Jungwook Choi “RA-LoRA: Rank-Adaptive Parameter-Efficient Fine-Tuning for Accurate 2-bit Quantized Large Language Models”, *In Findings of the Association for Computational Linguistics: ACL 2024*
- **[ACL 2024]** Janghwan Lee*, Seongmin Park*, Sukjin Hong, **Minsoo Kim**, Du-Seong Chang, and Jungwook Choi “Improving Conversational Abilities of Quantized Large Language Models via Direct Preference Alignment”, *In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*
- **[NeurIPS 2023]** Minsoo Kim, Sihwa Lee, Janghwan Lee, Sukjin Hong, Du-Seong Chang, Wonyong Sung and Jungwook Choi “Token-Scaled Logit Distillation for Ternary Weight Generative Language Models”, *Thirty-seventh Conference on Neural Information Processing Systems*. [Paper, Code]
- **[EMNLP 2023]** Janghwan Lee*, **Minsoo Kim***, Seungcheol Baek, Seok Joong Hwang, Wonyong Sung and Jungwook Choi “Enhancing Computation Efficiency in Large Language Models through Weight and Activation Quantization”, *In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics*. (*Co-First author) [Paper]
- **[EACL 2023]** Minsoo Kim, Kyuhong Shim, Seongmin Park, Wonyong Sung and Jungwook Choi, “Teacher Intervention: Improving Convergence of Quantization Aware Training for Ultra-Low Precision Transformers”, *In Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, pages 916–929, Dubrovnik, Croatia. Association for Computational*. [Paper, Code]
- **[EMNLP 2022]** Minsoo Kim, Sihwa Lee, Sukjin Hong, Du-Seong Chang, and Jungwook Choi, “Understanding and Improving Knowledge Distillation for Quantization-Aware Training of Large Transformer Encoders,” *In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 6713–6725, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics*. [Paper, Code]
- **[DAC 2022]** Joonsang Yu, Junki Park, Seongmin Park, **Minsoo Kim**, Sihwa Lee, Donghyun Lee, Jungwook Choi, “NN-LUT: neural approximation of non-linear operations for efficient transformer inference”, *In Proceedings of the 59th ACM/IEEE Design Automation Conference*. [Paper]

RESEARCH EXPERIENCE

- **Rank-Adaptive PEFT for 2-bit Quantized LLM Fine-Tuning (ACL 24)**
 - Identify inherent high-rank property of low-bit quantization error with thorough analysis
 - Investigate evolution of intrinsic subspace update in quantization combined LoRA fine-tuning
 - Propose rank adjusting method providing superior accuracy to SoTA quantized PEFT methods
- **Token-Scaled Logit Distillation for Ternary Weight Generative Language Models (NeurIPS 23)**
 - Investigate the challenges of applying Quantization-Aware Training (QAT) on a generative language model
 - Identify and analyze cumulative quantization error observed in causal attention of decoder model
 - Present confidence-based probabilistic correlation in the language modeling objective training
 - Propose novel KD method designed for GLM (up to 7B-sized) QAT, providing superior learning from teacher
- **Enhancing Efficiency in LLMs via Weight and Activation Quantization (EMNLP 23)**
 - Analyze various LLM (OPT, LLaMA) characteristics of weight/activation distribution with quantization
 - Scaling & calibration PTQ method effectively addressing combined weight and activation quantization effects
 - Identify underflow in W4A8; propose hybrid data format and arithmetic unit with $2\times$ HW efficiency
- **Understanding and Improving KD for QAT of Large Transformer Encoders (EMNLP 22)**
 - Mechanism of KD conducting attention recovery of quantized large Transformer encoders
 - Analyze quantization effect on attention behavior in Transformer over various target NLU tasks
 - Improve accuracy drop in NLU for 2bit weight quantization for large Transformer encoder with $<1\%$
- **Improving Convergence of QAT for Ultra-Low Precision Transformer Encoders (EACL 23)**
 - Proactive Teacher Intervention method for fast converging 2-bit QAT of Transformer encoders
 - Gradual Intervention Mechanism stabilizing the recovery of subsections of quantized Transformer layers
 - Achieve higher accuracy in BERT and ViT within up to 12.5x shorter fine-tuning time

SCHOLARSHIP AND AWARD

- **AICAS Grand Challenge 2024**, SW&HW Co-Optimization for LLM, 3rd place March 2024
- **Qualcomm Innovation Fellowship Korea 2023**, Winner, USD 3K, *Qualcomm* November 2023
- **Research Scholarship** USD 16K, *IoT System Semiconductor Research Center* Spring 2021 - Spring 2023
- **AI Grand Challenge**, *Korea Ministry of Science and ICT* November 2020
 - First place award in Model Compression Track (YOLOV5s Object Detection model 4x speed up)

SKILLS

- **Programming Languages:** Python, C, C++
- **Teaching Assistant:** SOC design (Spring 2021), Introduction to SW Optimization (Fall 2023)
- **English:** Served as a KATUSA (Korean Augmentation to the US Army) and certified Air Traffic Control Operator of the 8th Army (Jul 2017 - Apr 2019)
- **Academic Services:** Reviewer of EMNLP 23, ICML 24, ACL Rolling Review 24, COLM 24