BOYUAN ZHANG

+1 (509) 715-8390 \diamond bozhan@iu.edu \diamond Bloomington, IN

EDUCATION

Indiana University Bloomington	Aug. 2021 - Present
PhD in Intelligent System Engineering	
University of Southern California	Jan. 2019 - Dec. 2020
Master of Science in Electrical Engineering	
Shanghai Jiao Tong University	Sep. 2014 - Jun. 2018
Bachelor of Engineering in Information Engineering	

EXPERIENCE

Pacific Northwest National Laboratory

PhD Intern in High-Performance Computing

- Conducted an evaluation of lossy compression methods applied to microscopy images to assess their efficiency and effectiveness.
- Developed a new workflow on GPUs for AI-based compression techniques designed to achieve a high compression ratio, maintain quality, and ensure optimal performance.

PROJECTS

DLRM Communication Optimizations with Compression

Sep. 2022 - Feb. 2024

- Published in SC '24. Accelerating Communication in Deep Learning Recommendation Model Training with Dual-Level Adaptive Lossy Compression
- Developed multiple efficient compression schemes tailored to the specific data characteristics in DLRM, including a vector-based GPU LZ algorithm for embedding tables, to achieve a high compression ratio and better performance. The work also presents a dynamic error control scheme to manage error propagation, maintaining high accuracy, and a selection strategy for the best compression ratio among multiple schemes.

High-Performance AI-based Compression on GPUs Oct. 2023 - Jul. 2024

- Extended research from the PNNL internship. With a series of GPU optimizations, this work achieves significantly higher image quality at similar compression ratios compared to non-AI-based compressors and achieves even higher end-to-end performance.
- Utilized kernel fusion, warp-level optimization, shared memory, Prefix-Sum via Decoupled Lookback, GPU direct storage, and GPU pipeline techniques to achieve high performance.

High-Efficiency State Vector Quantum Circuit Simulation Mar. 2023 - Jan. 2024

• Utilized a unique characteristic of state vector simulation to divide the simulation process into separate jobs by partitioning the circuit. Leveraged compression to mitigate memory limitations in modern state vector quantum simulation algorithms, portable to simulators such as Qiskit-Aer, SV-Sim, cuQuantum, and Cirq. Achieved similar performance to state-ofthe-art simulators with significantly less memory consumption.

May 2023 - Sep. 2023

• Employed OpenMP to manage multiple GPUs and GPU streams, implementing an efficient pipeline. Also, proposed the first GPU-based point-wise relative error control scheme to limit the error produced by compression.

Fast and High-Ratio Lossy Compressor on GPUs Oct. 2022 - Jan. 2023

- Published in HPDC '23. FZ-GPU: A Fast and High-Ratio Lossy Compressor for Scientific Computing Applications on GPUs.
- Proposed a new pipeline to achieve both a high compression ratio and throughput by optimizing the dual-quantization in cuSZ to fit the new compression pipeline with significantly higher throughput, a new GPU bitshuffle process to fully leverage GPU parallelism, and a new fast lossless encoder to reduce redundancy introduced by bitshuffle.
- Utilized shared memory and warp-level vote function to implement an efficient bit-level memoryintensive operation, carefully managing shared memory to avoid bank conflicts, with kernel fusion for multiple processes.

Optimizing LZSS Lossless Compressor on GPUs Aug. 2022 - Jan. 2023

- Published in ICS '23. GPULZ: Optimizing LZSS Lossless Compression for Multi-byte Data on Modern GPUs
- Improved the state-of-the-art lossless LZ compressor (i.e., CULZSS) with significantly higher compression throughput. Designed a fully GPU-implemented LZ compressor that exploits the parallelism in the LZSS algorithm. Unlike state-of-the-art implementations, this work explores the advantages of utilizing multi-byte symbols in the sliding window lookup of the LZSS algorithm to achieve both higher compression ratio and performance.

SKILLS

CUDA, C++, C, Python

PUBLICATIONS

[SC '24] Hao Feng^{*}, Boyuan Zhang^{*}, Fanjiang Ye, Min Si, Ching-Hsiang Chu, Jiannan Tian, Chunxing Yin, Summer Deng, Yuchen Hao, Pavan Balaji, Tong Geng, Dingwen Tao. "Accelerating Communication in Deep Learning Recommendation Model Training with Dual-Level Adaptive Lossy Compression." The International Conference for High-Performance Computing, Networking, Storage, and Analysis, Atlanta, Georgia, USA, November 17-22, 2024. (*Equal contribution.)

[HPDC '23] Boyuan Zhang, Jiannan Tian, Sheng Di, Xiaodong Yu, Yunhe Feng, Xin Liang, Dingwen Tao, Franck Cappello. "FZGPU: A Fast and High-Ratio Lossy Compressor for Scientific Computing Applications on GPUs." The 32nd ACM International Symposium on High-Performance Parallel and Distributed Computing, Orlando, FL, June 16–23, 2023. DOI: 10.1145/3588195.3592994.

[ICS '23] Boyuan Zhang, Jiannan Tian, Sheng Di, Xiaodong Yu, Martin Swany, Dingwen Tao, Franck Cappello. "GPULZ: Optimizing LZSS Lossless Compression for Multi-byte Data on Modern GPUs." The 37th ACM International Conference on Supercomputing, Orlando, FL, USA, June 21–23, 2023. DOI: 10.1145/3577193.3593706.