

Xuhao Luo

9500 Gilman Drive, La Jolla, CA 92093
(858) 295-9996 \diamond x3luo@eng.ucsd.edu

Education

University of California San Diego M.S. in Computer Science, Department of Computer Science and Engineering GPA: 3.84/4.0	Sep 2019 - Mar 2021
University of Science and Technology of China (USTC) B.S. in Applied Physics, School of Physical Sciences Major in Microelectronics and Solid State Electronics	Sep 2015 - Jun 2019

Research Interests

Operating System, Networking, Computer Architecture, Heterogeneous Computing

Research Experience

An Asynchronous Executor for Distributed ML System <i>Research Project at Microsoft Research Asia</i>	Jun 2020 - Sep 2020
---	---------------------

- Designed and implemented an asynchronous executor for task scheduling and dispatching on multiple hardware.
- Designed and implemented CUDA-based high-performance inter-GPU communication channel for distributed ML within large-scale GPU cluster.
- Multi-GPU collective operation(AllReduce, AllGather, Broadcast) throughput outperforms Nvidia NCCL by at most 18.4% under the same system setting.

An FPGA-based Disaggregated Memory System <i>Research Project at UCSD</i>	Sep 2019 - Jun 2020 <i>Supervisor: Prof. Yiyang Zhang</i>
---	--

- Working on FPGA-based disaggregated virtual memory system for system resource disaggregation.
- Designed and implemented a go-back-N based full reliable network stack on both FPGA and host Linux server to support high-performance reliable network communication. Using kernel-bypass to achieve high-throughput and low-latency.
- Designed and implemented an RPC-semantic connectionless network stack to improve scalability, with a delay-based congestion control.
- Achieved RDMA-like latency and 10Gbps(limited by hardware interface) throughput at rack-scale.

Past Projects

Design and Implementation of HLS Based Quantized Neural Network Accelerator <i>Graduation Project</i>	Jan 2019 - May 2019 <i>Supervisor: Prof. Xi Jin</i>
---	--

- Studied the 8-bit quantization algorithm, including the quantization algorithm, the dequantization algorithm and the implementation of the 8-bit quantized convolution.
- Designed and implemented a general 8-bit quantized convolution module on Xilinx Virtex FPGA, which achieved high parallelization through array architecture, and realized memory access optimization through data reuse.
- Developed the TensorFlow C++ API for the hardware accelerator using OpenCL. Used this accelerator to accelerate the ResNet-50 CNN and achieved a speedup of 5.17x and a memory usage reduction of 66% compared with the CPU TensorFlow implementation on Xeon E5 2686.

Binary Neuron Network (BNN) Acceleration using HLS <i>Summer Internship at Cornell University</i>	Jul 2018 - Sep 2018 <i>Supervisor: Prof. Zhiru Zhang</i>
---	---

- Designed and implemented a BNN accelerator for LeNet-5 for MNIST handwritten digits recognition.

- Applied multiple methods to improve the performance of the accelerator including parallelization, pipelining, line buffer, task-level parallelism and batch processing.
- Implemented the accelerator on Zedboard, ZC706, and AWS EC2 F1. Achieved speedups of 33x(580fps), 88x(1543fps) and 114x(2170fps) compared with the software implementation baseline on Intel Xeon 5420 CPU.

Honors and Awards

- | | |
|---|----------|
| • USTC Class of 2019 Outstanding Graduates | May 2019 |
| • 2017/18 USTC Outstanding Students Scholarship, Golden Award | Sep 2018 |
| • 2016/17 USTC Outstanding Students Scholarship, Silver Award | Sep 2017 |
| • 2015/16 USTC Outstanding Students Scholarship, Bronze Award | Sep 2016 |
| • The 13 th Competition of Physical Research Experiment, 2 nd Prize | Dec 2017 |
| • The 6 th Aegon-Industrial Fund Scholarship | Jun 2017 |

Skills

- | | |
|------------------------|---|
| Language | C/C++, Python, Go, Rust, Haskell, OpenCL, Verilog, HTML, JavaScript |
| Tools/Framework | TensorFlow, Docker, Zookeeper, LLVM, Google Test |