

JIECAO YU

CONTACT INFORMATION

Website: <https://jiecaoyu.github.io/>
E-mail: jiecaoyu@fb.com
Tel: +1 (734) 353 - 8285

RESEARCH INTERESTS

Software-hardware co-design for deep learning acceleration, DNN pruning and quantization.

EDUCATION

Ph.D. Candidate, Computer Science & Engineering 08/2014-09/2019

Advisor: Prof. Scott Mahlke

University of Michigan, Ann Arbor, MI

Dissertation: Efficient Deep Neural Network Computation on Processors

M.S. Computer Science & Engineering 08/2014-12/2015

University of Michigan, Ann Arbor, MI

Cumulative GPA: 4.00/4

B.Eng. Electronic & Information Engineering 08/2010-06/2014

Honored Minor, Advanced Honor Class of Engineering Education (ACEE)

Zhejiang University, Hangzhou, China

Cumulative GPA: 92/100 (3.98/4.0), Rank: 2/92

PUBLICATIONS

Mao Ye, Dhruv Choudhary, Jiecao Yu, Ellie Wen, Zeliang Chen, Jiyan Yang, Jongsoo Park, Qiang Liu, Arun Kejariwal. “*Adaptive Dense-to-Sparse Paradigm for Pruning Online Recommendation System with Non-Stationary Data*”. preprint at arXiv: 2010.08655

Jiecao Yu, Andrew Lukefahr, Reetuparna Das, Scott Mahlke. “*TF-Net: Deploying Sub-Byte Deep Neural Networks on Microcontrollers*”. ESWEK-TECS special issue / the International Conference on Compilers, Architecture, and Synthesis for Embedded Systems (CASES), Oct, 2019

Jiecao Yu, Jongsoo Park, Maxim Naumov. “*Spatial-Winograd Pruning Enabling Sparse Winograd Convolution*”. preprint at arXiv: 1901.02132

Xiaowei Wang, Jiecao Yu, Charles Augustine, Ravi Iyer, Reetuparna Das. “*Bit Prudent In-Cache Acceleration of Deep Convolutional Neural Networks*”. The 25th International Symposium on High-Performance Computer Architecture (HPCA-25), Feb, 2019

Jiecao Yu, Andrew Lukefahr, David Palframa, Ganesh Dasika, Reetuparna Das, Scott Mahlke. “*Scalpel: Customizing DNN Pruning to the Underlying Hardware Parallelism*”. The 44th International Symposium on Computer Architecture (ISCA-44), Jun, 2017

Jiecao Yu, Andrew Lukefahr, Shruti Padmanabha, Reetuparna Das, Scott Mahlke. “*Adaptive Cache Partitioning on a Composite Core*”. The PRISM-3 Workshop at the International Symposium on Computer Architecture (ISCA-42), Jun, 2015

EXPERIENCES

Facebook, Inc.

Research Scientist, AI System SW/HW Co-design Group
Manager: Dr. Jongsoo Park

10/2019-Present
Menlo Park, CA

- Working on DNN model pruning and acceleration.

University of Michigan

Graduate Student Research Assistant

08/2014-09/2019
Ann Arbor, MI

- Investigating the training algorithms of binary/ ternary neural networks.
- Developing low-precision computation algorithms/ hardware architecture for mobile and embedded devices.
- Developed a new DNN pruning technique, Scalpel, which applies weight pruning and node pruning synergistically based on the underlying hardware platform to improve the computation performance.

Facebook, Inc.

Research Intern, AI System SW/HW Co-design Group
Manager: Dr. Jongsoo Park

05/2018-08/2018
Menlo Park, CA

- Proposed a two-step pruning technique, spatial-Wingorad pruning, to improve the Winograd-domain sparsity.

Arm, Inc.

Research Intern, Machine Learning Group
Manager: Dr. Ganesh Dasika

05/2017-07/2017
Austin, TX

- Profiling and analysis of image captioning workloads (Show-and-Tell/Show-Attend-and-Tell).
- Built and profiled the server-side image captioning/classifying workloads based on TensorFlow Serving.

Arm, Inc.

Research Intern, Machine Learning Group
Manager: Dr. David Palframan

06/2016-08/2016
Austin, TX

- Worked on Deep Neural Network acceleration on Arm cores, especially low-power microcontrollers.
- DNN weight pruning techniques are employed to compress the DNN in the keyword spotting (KWS) system.
- Libraries for sparse matrix computation on Arm Cortex-M4 microcontrollers are implemented and well-optimized.

University of Southern California

Research Intern
Supervisor: Prof. Melvin Breuer

07/2013-09/2013
Los Angeles, CA

- Worked on enhancing yield of VLSI chips via redundancy.

PATENTS

US 20180373975, "Systems and Devices for Compressing Neural Network Parameters", Jiecao Yu, Andrew Lukefahr, David Palframan, Ganesh Dasika, Reetuparnda Das, Scott Mahlke, Filed: June 21, 2017

US 20180373978, "Systems and Devices for Formatting Neural Network Parameters", Jiecao Yu, Andrew Lukefahr, David Palframan, Ganesh Dasika, Reetuparnda Das, Scott Mahlke, Filed: June 21, 2017

US 20170262285, "Controlling Transition Between Using First and Second Processing Circuitry", Andrew Lukefahr, Shruti Padmanabha, Jiecao Yu, Reetuparna Das, and Scott Mahlke, Filed: March 08, 2016

| | |
|------------------------------------|---|
| TALKS & POSTERS | <p>[Talk] Jiecao Yu. “<i>Efficient Low-Precision Deep Neural Networks on IoT Microcontrollers</i>”. Arm Research Summit, Sep, 2019</p> <p>[Poster] Babak Zamirai, Jiecao Yu, Salar Latifi, Scott Mahlke. “<i>Input-specialized Heterogeneous Neural Networks</i>”. C-FAR 2016 Annual Meeting, Dec, 2016</p> <p>[Poster] Salar Latifi, Babak Zamirai, Jiecao Yu, Scott Mahlke. “<i>Quality Assurance for Approximate Computing</i>”. C-FAR 2016 Annual Meeting, Dec, 2016</p> <p>[Poster] Jiecao Yu, Babak Zamirai, Scott Mahlke. “<i>An Interactive Deep Neural Network Pruning System</i>”. C-FAR 2016 Semi-Annual Meeting, May, 2016</p> |
| SERVICE | <p>Reviewer:</p> <ul style="list-style-type: none"> • Design Automation Conference (DAC) Technical Program Committee (TPC) Member ('20) • IEEE Transactions on Computers ('20) • IEEE Access ('19, 20) • IEEE Transactions on Neural Networks and Learning Systems (TNNLS'19, 20) • Elsevier Journal of Systems Architecture (JSA'19) • ACM Journal on Emerging Technologies in Computing Systems (JETC'17) <p>Second Reviewer:</p> <ul style="list-style-type: none"> • Int. Conference on Compilers, Architecture, and Synthesis for Embedded Systems (CASES'17, 18, 19) • Int. Symposium on Microarchitecture (MICRO'17, 19) • Int. Symposium on Computer Architecture (ISCA'15, 17) • Int. Symposium on Code Generation and Optimization (CGO'16, 17) • Int. Symposium on High-Performance Computer Architecture (HPCA'16, 17) • Int. Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'17) • Int. Conference on Supercomputing (ICS'16) |
| COURSE PROJECTS | <p>L1 Cache Partitioning on a SMT Core Winter 2015 <i>Parallel Computer Architecture (EECS 570), Prof. Thomas Wenisch</i> Designed a dynamic L1 cache partitioning mechanism for a SMT core based on way-partitioning technique and augmented LRU replacement policy. Cache capacities can be resized at a fine granularity to capture the change of the cache demands for different threads.</p> <p>A Two-Way Superscalar R10K SMT Processor Fall 2014 <i>Computer Architecture (EECS 470), Prof. Trevor Mudge</i> Designed and implemented a synthesizable two-way superscalar Out-of-Order processor in Verilog HDL with speculative LSQ, instruction prefetching and supporting of simultaneous multithreading.</p> |
| RELEVANT GRADUATE COURSEWORK | <p>University of Michigan - Ann Arbor</p> <ul style="list-style-type: none"> • EECS 470: Computer Architecture (A+) • EECS 583: Advanced Compilers (A+) • EECS 570: Parallel Computer Architecture (A) • EECS 492: Introduction to Artificial Intelligence (A+) • EECS 573: Microarchitecture (A) |

| | | |
|------------------------|---|------------------|
| AWARDS & HONORS | National Scholarship (top 1.8%), <i>China</i> | 2011 |
| | First-Class Scholarship of National IC Talents Training Base, <i>China</i> | 2012, 2013 |
| | First-Class Scholarship for Outstanding Students (top 3%), <i>Zhejiang University</i> | 2011, 2012, 2013 |
| | Honorable Mention in MCM/ICM Contest, <i>United States</i> | 2013 |
| SKILLS | Language proficiency: Fluent English, Native Chinese | |
| | Programming: Python, C/C++, Bash, L ^A T _E X, Verilog HDL, VHDL, MATLAB | |
| | Tools: Caffe, Torch/PyTorch, TensorFlow, LLVM, Gem5 | |
| TEACHING EXPERIENCE | CMOS Integrated Circuits Design | Fall 2013 |
| | College of Electrical Engineering | |
| | Zhejiang University, Hangzhou, China | |