

2019 ACM/IEEE 46th Annual International Symposium on Computer Architecture (ISCA 2019)

**Phoenix, Arizona, USA
22 – 26 June 2019**



**IEEE Catalog Number: CFP19030-POD
ISBN: 978-1-7281-4838-0**

**Copyright © 2019, The Association for Computing Machinery (ACM)
All Rights Reserved**

****** This is a print representation of what appears in the IEEE Digital Library. Some format issues inherent in the e-media version may also appear in this print version.***

IEEE Catalog Number:	CFP19030-POD
ISBN (Print-On-Demand):	978-1-7281-4838-0
ISBN (Online):	978-1-4503-6669-4
ISSN:	1063-6897

Additional Copies of This Publication Are Available From:

Curran Associates, Inc
57 Morehouse Lane
Red Hook, NY 12571 USA
Phone: (845) 758-0400
Fax: (845) 758-2633
E-mail: curran@proceedings.com
Web: www.proceedings.com

CURRAN ASSOCIATES INC.
proceedings
.com

Contents

ISCA-46 Sponsors	vii
General Chair Foreword	x
Program Chair Foreword	xii
ISCA-46 Organizing Committee	xxiii
ISCA-46 Program Committee	xxiv
ISCA-46 External Review Committee	xxvi
Perceptron-Based Prefetch Filtering	1
Eshan Bhatia (<i>Texas A&M University</i>); Daniel A. Jiménez (<i>Barcelona Supercomputing Center / Texas A&M University</i>); Paul Gratz (<i>Texas A&M University</i>); Elvira Teran (<i>Texas A&M International University</i>); Seth Pugsley (<i>Intel Labs</i>); Gino Chacon (<i>Texas A&M University</i>)	
Post-Silicon CPU Adaptations Made Practical Using Machine Learning	14
Stephen J Tarsa (<i>Intel Labs</i>); Rangeen Basu Roy Chowdhury, Julien Sebot (<i>Intel Corporation</i>); Gautham Chinya (<i>Intel Labs</i>); Jayesh Gaur (<i>Intel Labs India</i>); Karthik Sankaranarayanan, Chit-Kwan Lin (<i>Intel Labs</i>); Robert Chappell, Ronak Singhal (<i>Intel Corporation</i>); Hong Wang (<i>Intel Labs</i>)	
Bit-Level Perceptron Prediction for Indirect Branch Prediction	27
Elba Garza, Samira Mirbagher Ajorpaz, Tahsin Ahmad Khan (<i>Texas A&M University</i>); Daniel A. Jiménez (<i>Barcelona Supercomputing Center and Texas A&M University</i>)	
Generative and Multi-phase Learning for Computer Systems Optimization	39
Yi Ding, Nikita Mishra, Henry Hoffmann (<i>University of Chicago</i>)	
OO-VR: NUMA Friendly Object-Oriented VR Rendering Framework For Future NUMA-Based Multi-GPU Systems	53
Chenhao Xie (<i>University of Houston; Pacific Northwest National Laboratory</i>); Xin Fu (<i>University of Houston</i>); Mingsong Chen (<i>East China Normal University</i>); Shuaiwen Leon Song (<i>Pacific Northwest National Laboratory; The University of Sydney</i>)	
PES: Proactive Event Scheduling for Energy-Efficient Mobile Web Computing	66
Yu Feng, Yuhao Zhu (<i>University of Rochester</i>)	
3D-based Video Understanding Acceleration by Leveraging Temporal Locality and Activation Sparsity	79
Huixiang Chen, Mingcong Song, Jiechen Zhao (<i>University of Florida</i>); Yuting Dai (<i>Guizhou University</i>); Tao Li (<i>University of Florida</i>)	
Energy-Efficient Video Processing for Virtual Reality	91
Yue Leng (<i>UIUC</i>); Chi-chun Chen, Qiuyue Sun (<i>University of Rochester</i>); Jian Huang (<i>UIUC</i>); Yuhao Zhu (<i>University of Rochester</i>)	
Triad-NVM: Persistency for Integrity-Protected and Encrypted Non-Volatile Memories	104
Amro Awad, Yan Solihin (<i>University of Central Florida</i>); Laurent Njilla (<i>Air Force Research Lab</i>); Mao Ye, Kazi Zubair (<i>University of Central Florida</i>)	

GraphSSD: Graph Semantics Aware SSD	116
Kiran kumar matam (<i>USC</i>); Gunjae Koo (<i>Hongik University</i>); Haipeng Zha (<i>USC</i>); Hung-Wei Tseng (<i>NCSU</i>); Murali Anavarum (<i>USC</i>)	
CROW: A Low-Cost Substrate for Improving DRAM Performance, Energy Efficiency, and Reliability	129
Hasan Hassan, Minesh Patel (<i>ETH Zurich</i>); Jeremie S. Kim (<i>Carnegie Mellon University; ETH Zurich</i>); A. Giray Yaglikci (<i>ETH Zurich</i>); Nandita Vijaykumar (<i>Carnegie Mellon University; ETH Zurich</i>); Nika Mansouri Ghiasi (<i>ETH Zurich</i>); Saugata Ghose (<i>Carnegie Mellon University</i>); Onur Mutlu (<i>ETH Zurich; Carnegie Mellon University</i>)	
Janus: Optimizing Memory and Storage Support for Non-Volatile Memory Systems	143
Sihang Liu, Korakit Seemakhupt (<i>University of Virginia</i>); Gennady Pekhimenko (<i>University of Toronto</i>); Aasheesh Kolli (<i>Pennsylvania State University and VMware Research</i>); Samira Khan (<i>University of Virginia</i>)	
Anubis: Low-Overhead and Practical Recovery Time for Secure Non-Volatile Memories	157
Kazi Abu Zubair, Amro Awad (<i>University of Central Florida</i>)	
Emerald: Graphics Modeling for SoC Systems	169
Ayub A. Gubran, Tor M. Aamodt (<i>University of British Columbia</i>)	
Linebacker: Preserving Victim Cache Lines in Idle Register Files of GPUs	183
Yunho Oh (<i>EPFL</i>); Gunjae Koo (<i>Hongik University</i>); Murali Annavaram (<i>University of Southern California</i>); Won Woo Ro (<i>Yonsei University</i>)	
MGPUSim: Enabling Multi-GPU Performance Modeling and Optimization	197
Yifan Sun, Trinayan Baruah (<i>Northeastern University</i>); Saiful A. Mojumder (<i>Boston University</i>); Shi Dong, Xiang Gong, Shane Treadway, Yuhui Bao, Spencer Hance, Carter McCardwell, Vincent Zhao, Harrison Barclay (<i>Northeastern University</i>); Amir Kavayan Ziabari, Zhongliang Chen (<i>AMD</i>); Rafael Ubal (<i>Northeastern University</i>); José L. Abellán (<i>Universidad Católica San Antonio de Murcia</i>); John Kim (<i>KAIST</i>); Ajay Joshi (<i>Boston University</i>); David Kaeli (<i>Northeastern University</i>)	
Opportunistic Computing in GPU Architectures	210
Ashutosh Pattnaik, Xulong Tang (<i>Penn State</i>); Onur Kayiran (<i>AMD Research</i>); Adwait Jog (<i>College of William & Mary</i>); Asit Mishra (<i>NVIDIA</i>); Mahmut T. Kandemir, Anand Sivasubramaniam, Chita Das (<i>Penn State</i>)	
Interplay between Hardware Prefetcher and Page Eviction Policy in CPU-GPU Unified Virtual Memory	224
Debashis Ganguly, Rami Melhem (<i>Department of Computer Science, University of Pittsburgh</i>); Jun Yang (<i>Electrical and Computer Engineering, University of Pittsburgh</i>); Ziyu Zhang (<i>Department of Computer Science, University of Pittsburgh</i>)	
Sparse ReRAM Engine: Joint Exploration of Activation and Weight Sparsity in Compressed Neural Networks	236
Tzu-Hsien Yang (<i>National Taiwan University</i>); Hsiang-Yun Cheng (<i>Academia Sinica</i>); Chia-Lin Yang, I-Ching Tseng (<i>National Taiwan University</i>); Han-Wen Hu, Hung-Sheng Chang, Hsiang-Pang Li (<i>Macronix International Co., Ltd.</i>)	
MnnFast: A Fast and Scalable System Architecture for Memory-Augmented Neural Networks . .	250
Hanhwi Jang (<i>POSTECH</i>); Joonsung Kim (<i>Seoul National University</i>); Jae-Eon Jo (<i>POSTECH</i>); Jaewon Lee, Jangwoo Kim (<i>Seoul National University</i>)	

TIE: Energy-efficient tensor train-based inference engine for deep neural network	264
Chunhua Deng (<i>Rutgers University</i>); Fangxuan Sun (<i>Nanjing University</i>); Xuehai Qian (<i>University of Southern California</i>); Jun Lin, Zhongfeng Wang (<i>Nanjing University</i>); Bo Yuan (<i>Rutgers University</i>)	
Accelerating Distributed Reinforcement Learning with In-Switch Computing	279
Youjie Li, Iou-Jen Liu, Yifan Yuan, Deming Chen, Alexander Schwing, Jian Huang (<i>UIUC</i>)	
Eager Pruning: Algorithm and Architecture Support for Fast Training of Deep Neural Networks .	292
Jiaqi Zhang, Xiangru Chen, Mingcong Song, Tao Li (<i>University of Florida</i>)	
Laconic Deep Learning Inference Acceleration	304
Sayeh Sharify, Alberto Delmas Lascorz, Mostafa Mahmoud, Milos Nikolic, Kevin Siu, Dylan Malone Stuart, Zissis Poulos, Andreas Moshovos (<i>Toronto</i>)	
MicroScope: Enabling Microarchitectural Replay Attacks	318
Dimitrios Skarlatos, Mengjia Yan, Bhargava Gopireddy, Read Sprabery, Josep Torrellas, Christopher Fletcher (<i>University of Illinois at Urbana Champaign</i>)	
SecDir: A Secure Directory to Defeat Directory Side-Channel Attacks	332
Mengjia Yan, Jen-Yang Wen, Christopher Fletcher, Josep Torrellas (<i>University of Illinois at Urbana Champaign</i>)	
Secure TLBs	346
Shuwen Deng, Wenjie Xiong, Jakub Szefer (<i>Yale University</i>)	
New Attacks and Defense for Encrypted-Address Cache	360
Moinuddin K Qureshi (<i>Georgia Tech</i>)	
InvisiPage: Oblivious Demand Paging for Secure Enclaves	372
Shaizeen Aga, Satish Narayanasamy (<i>University of Michigan, Ann Arbor</i>)	
TWiCe: Preventing Row-hammering by Exploiting Time Window Counters	385
Eojin Lee, Ingab Kang (<i>Seoul National University</i>); Sukhan Lee (<i>Seoul National University / Samsung Electronics</i>); G. Edward Suh (<i>Cornell University</i>); Jung Ho Ahn (<i>Seoul National University</i>)	
Duality Cache for Data Parallel Acceleration	397
Daichi Fujiki, Scott Mahlke, Reetuparna Das (<i>University of Michigan</i>)	
Adaptive Memory-Side Last-Level GPU Caching	411
Xia Zhao, Almutaz Adileh (<i>Ghent University</i>); Zhibin Yu (<i>Shenzhen Institutes of Advanced Technology, CAS</i>); Zhiying Wang (<i>National University of Defense Technology</i>); Aamer Jaleel (<i>NVIDIA</i>); Lieven Eeckhout (<i>Ghent University</i>)	
SCU: A GPU Stream Compaction Unit for Graph Processing	424
Albert Segura, Jose-Maria Arnau, Antonio Gonzalez (<i>Universitat Politecnica de Catalunya</i>)	
Filter Caching for Free: The Untapped Potential of the Store Buffer	436
Ricardo Alves (<i>Uppsala University</i>); Alberto Ros (<i>Universidad de Murcia</i>); David Black-Schaffer, Stefanos Kaxiras (<i>Uppsala University</i>)	
Efficient Meta-Data Management for Irregular Data Prefetching	449
Hao Wu (<i>ARM, UT Austin</i>); Krishnendra Nathella, Dam Sunwoo (<i>ARM</i>); Akanksha Jain, Calvin Lin (<i>UT Austin</i>)	

AsmDB: Understanding and Mitigating Front-end Stalls in Warehouse-Scale Computers	462
Grant Ayers (<i>Stanford University</i>); Nayana Prasad Nagendra, David I. August (<i>Princeton University</i>); Hyoun Kyu Cho, Svilen Kanev (<i>Google</i>); Christos Kozyrakis (<i>Stanford University</i>); Trivikram Krishnamurthy (<i>Nvidia</i>); Heiner Litz (<i>UC Santa Cruz</i>); Tipp Moseley, Parthasarathy Ranganathan (<i>Google</i>)	
Fine-grained Warm Water Cooling for Improving Datacenter Economy	474
Weixiang Jiang, Ziyang Jia, Sirui Feng, Fangming Liu, Hai Jin (<i>Huazhong University of Science and Technology</i>)	
DeepAttest: An End-to-End Attestation Framework for Deep Neural Networks	487
Huili Chen, Cheng Fu, Bitar Darvish Rouhani, Jishen Zhao, Farinaz Koushanfar (<i>UC San Diego</i>)	
A Time-Space Sharing Scheduling Abstraction for Next Generation of Shared Cloud via Vertical Labels	499
Yuzhao Wang (<i>Huazhong University of Science and Technology, Shenzhen Institute of Advanced Technology</i>); Lele Li (<i>Shenzhen Institute of Advanced Technology</i>); You Wu (<i>University of Southern California</i>); Junqing Yu (<i>Huazhong University of Science and Technology</i>); Zhibin Yu (<i>Shenzhen Institute of Advanced Technology</i>); Xuehai Qian (<i>University of Southern California</i>)	
SoftSKU: Optimizing Server Architectures for Microservice Diversity @Scale	513
Akshitha Sriraman (<i>University of Michigan</i>); Abhishek Dhanotia (<i>Facebook</i>); Thomas F. Wenisch (<i>University of Michigan</i>)	
Full-Stack, Real-System Quantum Computer Studies: Architectural Comparisons and Design Insights	527
Prakash Murali (<i>Princeton University</i>); Norbert Matthias Linke (<i>University of Maryland</i>); Margaret Martonosi (<i>Princeton University</i>); Ali Javadi Abhari (<i>IBM T. J. Watson Research Center</i>); Nhung Hong Nguyen, Cinthia Huerta Alderete (<i>University of Maryland</i>)	
Statistical Assertions for Validating Patterns and Finding Bugs in Quantum Programs	541
Yipeng Huang, Margaret Martonosi (<i>Princeton University</i>)	
Asymptotic Improvements to Quantum Circuits via Qutrits	554
Pranav Gokhale, Jonathan M. Baker, Casey Duckering (<i>University of Chicago</i>); Natalie C. Brown (<i>Georgia Institute of Technology</i>); Kenneth R. Brown (<i>Duke University</i>); Frederic T. Chong (<i>University of Chicago</i>)	
A Stochastic-Computing based Deep Learning Framework using Adiabatic Quantum-Flux-Parametron Superconducting Technology	567
Ruizhe Cai, Ao Ren (<i>Northeastern University</i>); Olivia Chen (<i>Yokohama National University</i>); Ning Liu, Caiwen Ding (<i>Northeastern University</i>); Xuehai Qian (<i>University of Southern California</i>); Jie Han (<i>University of Alberta</i>); Wenhui Luo, Nobuyuki Yoshikawa (<i>Yokohama National University</i>); Yanzhi Wang (<i>Northeastern University</i>)	
A Quantum Computational Compiler and Design Tool for Technology-Specific Targets	579
Kaitlin Smith, Mitch Thornton (<i>Southern Methodist University</i>)	
IntelliNoC: A Holistic Design Framework for Energy-Efficient and Reliable On-Chip Communication for Manycores	589
Ke Wang, Ahmed Louri (<i>George Washington University</i>); Avinash Karanth, Razvan Bunescu (<i>Ohio University</i>)	

HALO: Accelerating Flow Classification for Scalable Packet Processing in NFV	601
Yifan Yuan (<i>UIUC</i>); Yipeng Wang, Ren Wang (<i>Intel Labs</i>); Jian Huang (<i>UIUC</i>)	
Scalable Interconnects for Reconfigurable Spatial Architectures	615
Yaqi Zhang, Alexander Rucker, Matthew Viliam, Raghu Prabhakar, William Hwang, Kunle Olukotun (<i>Stanford University</i>)	
CoNDA: Efficient Cache Coherence Support for Near-Data Accelerators	629
Amirali Boroumand, Saugata Ghose (<i>Carnegie Mellon University</i>); Minesh Patel (<i>Carnegie Mellon University, ETH</i>); Rachata Ausavarungnirun (<i>Carnegie Mellon University, King Mongkut's University of Technology North Bangkok</i>); Hasan Hassan (<i>Carnegie Mellon University, ETH</i>); Brandon Lucia, Kevin Hsieh (<i>Carnegie Mellon University</i>); Nastaran Hajinazar (<i>Carnegie Mellon University, Simon Fraser University</i>); Krishna T. Malladi, Hongzhong Zheng (<i>Samsung Semiconductor, Inc.</i>); Onur Mutlu (<i>ETH, Carnegie Mellon University</i>)	
Designing Vertical Processors in Monolithic 3D	643
Bhargava Gopireddy, Josep Torrellas (<i>University of Illinois at Urbana Champaign</i>)	
Time Squeezing for Tiny Devices	657
Yuanbo Fan, Simone Campanoni, Russ Joseph (<i>Northwestern University</i>)	
XPC: Architectural Support for Secure and Efficient Cross Process Call	671
Dong Du, Zhichao Hua, Yubin Xia, Binyu Zang, Haibo Chen (<i>Shanghai Jiao Tong University</i>)	
AxMemo: Hardware-Compiler Co-Design for Approximate Code Memoization	685
Zhenhong Liu (<i>UIUC</i>); Amir Yazdanbakhsh (<i>Google Brain</i>); Dong Kai Wang (<i>UIUC</i>); Hadi Esmaeilzadeh (<i>UCSD</i>); Nam Sung Kim (<i>UIUC</i>)	
Translation Ranger: Operating System Support for Contiguity-Aware TLBs	698
Zi Yan (<i>Rutgers University/NVIDIA</i>); Daniel Lustig, David Nellans (<i>NVIDIA</i>); Abhishek Bhattacharjee (<i>Yale University</i>)	
Bouncer: Static Program Analysis in Hardware	711
Joseph McMahan, Michael Christensen (<i>University of California, Santa Barbara</i>); Kyle Dewey (<i>California State University Northridge</i>); Ben Hardekopf, Timothy Sherwood (<i>University of California, Santa Barbara</i>)	
Efficient Invisible Speculative Execution through Selective Delay and Value Prediction	723
Christos Sakalis, Stefanos Kaxiras (<i>Uppsala University</i>); Alberto Ros (<i>University of Murcia</i>); Alexandra Jimborean (<i>Uppsala University</i>); Magnus Sjölander (<i>Norwegian University of Science and Technology</i>)	
Stream-based Memory Access Specialization for General Purpose Processors	736
Zhengrong Wang, Tony Nowatzki (<i>UCLA</i>)	
Using SMT to accelerate nested virtualization	750
Lluís Vilanova (<i>Technion</i>); Nadav Amit (<i>VMWare Research</i>); Yoav Etsion (<i>Technion</i>)	
Master of None Acceleration: A Comparison of Accelerator Architectures for Analytical Query Processing	762
Andrea Lottarini (<i>Google</i>); João P. Cerqueira, Thomas J. Repetti, Stephen A. Edwards, Kenneth A. Ross, Mingoo Seok, Martha A. Kim (<i>Columbia University</i>)	

Cryogenic Computer Architecture Modeling with Memory-Side Case Studies	774
Gyu-Hyeon Lee, Dongmoon Min, Il-Kwon Byun, Jangwoo Kim (<i>Department of Electrical and Computer Engineering, Seoul National University</i>)	
Cambricon-F: Machine Learning Computers with Fractal von Neumann Architecture	788
Yongwei Zhao, Zidong Du, Qi Guo, Shaoli Liu (<i>Institute of Computing Technology, Chinese Academy of Sciences</i>); Ling Li (<i>Institute of Software, Chinese Academy of Sciences</i>); Zhiwei Xu, Tianshi Chen, Yunji Chen (<i>Institute of Computing Technology, Chinese Academy of Sciences</i>)	
FloatPIM: In-Memory Acceleration of Deep Neural Network Training with High Precision	802
Mohsen Imani, Saransh Gupta, Yeseong Kim, Tajana Rosing (<i>UC San Diego</i>)	
Author index	816
ISCA-46 Other External Reviewers	819