# 2023 56th IEEE/ACM International Symposium on Microarchitecture (MICRO 2023)

Toronto, Ontario, Canada
28 October - 1 November 2023

Pages 1-741

*** *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.*

**Additional Copies of This Publication Are Available From:**

# Table of contents

## Best Paper Session
Session Chair:  Davide Basilio Bartolini (Huawei)

## Session 1A: Accelerators Based on HW/SW Co-Design Accelerators for Matrix Processing

Session Chair:  Michael Pellauer (NVIDIA)

## Session 1B: Architectural Support/Programming Languages, Case Study

Session Chair:  Saugata Ghose (University of Illinois Urbana-Champaign)

## Session 1C: Design Automation, Synthesis, Hardware Generation

Session Chair: Mark Jeffrey (University of Toronto)

## Session 2A: ML Design Space Exploration Generation

Session Chair: Tushar Krishna (Georgia Institute of Technology)

## Session 2B: Microarchitecture

Session Chair: Daniel Sorin (Duke University)

## Session 2C: Accelerators for Graphs, Robotics

Session Chair: Sabrina Neuman (Boston University)

## Session 3A: ML Sparsity

Session Chair: Biswabandan Panda (Indian Institute of Technology Bombay)

## Session 3B: GPUs

Session Chair:  Nandita Vijaykumar (University of Toronto); Sabrina Neuman (Boston University)

## Session 4A: ML Architecture

Session Chair:  Po-An Tsai (NVIDIA)

## Session 4B: Quantum

Session Chair: Hiroaki Kobayashi (Tohoku University)

## Session 4C: Emerging Technologies: Superconducting, Photonics, DNA

Session Chair:  Koji Inoue (Kyushu University)

## Session 5A: Security Encryption/Confidentiality Support

Session Chair: Gururaj Saileshwar (University of Toronto / NVIDIA Research)

# Session 5B: Prefetching

Session Chair:  Leeor Peled (Toga Networks)

# Session 5C: Processing-In-Memory

Session Chair:  Dimitrios Skarlatos (Carnegie Mellon University)

# Session 6A: Security Hardware

Session Chair:  Samira Mirbagher Ajorpaz (North Carolina State University)

# Session 6B: Datacenter Networks

Session Chair:  Trevor E. Carlson (National University of Singapore)

# Session 6C: Reliability, Availability

Session Chair: Freddy Gabbay (Ruppin Academic College)

# Session 7A: Accelerators Various

Session Chair: Alex K. Jones (University of Pittsburgh)

## Session 7B: Caches, Intermitent Computing, Persistency

Session Chair:  Rachata Ausavarungnirun (King Mongkut's University of Technology North Bangkok)

## Session 8A: Accelerators for Neural Nets <br> Accelerators for Matrix Processing

Session Chair:  Jason Clemons (NVIDIA)

# Session 8B: Virtual Memory (Translation)

Session Chair: Mohammad Alian (University of Kansas)

# Session 8C: Benchmarking and Methodology

Session Chair: Miquel Moretó (Universitat Politècnica de Catalunya/Barcelona Supercomputing Center)

## Session 9A: Accelerators in Processors

Session Chair:  Sihang Liu (University of Waterloo)

## Session 9B: ML Compiler Optimizations/Reconfigurable Architectures

Session Chair:  Jian Huang (University of Illinois Urbana-Champaign)

## Session 9C: Domain Specific Genomics

Session Chair:  Pradip Bose (IBM)