# 2024 ACM/IEEE 51st Annual International Symposium on Computer Architecture (ISCA 2024)

Buenos Aires, Argentina
29 June - 3 July 2024

Pages 1-659

**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

# 2024 ACM/IEEE 51st Annual International Symposium on Computer Architecture (ISCA)

# ISCA 2024

## Table of Contents

## Session 1A: Microarchitectire

*Ishita Chaturvedi (Princeton University, USA), Bhargav Reddy Godala (Princeton University, USA), Yucan Wu (Princeton University, USA), Ziyang Xu (Princeton University, USA), Konstantinos Iliakis (National Technical University of Athens, Greece), Panagiotis-Eleftherios Eleftherakis (National Technical University of Athens, Greece), Sotirios Xydis (National Technical University of Athens, Greece), Dimitrios Soudris (National Technical University of Athens, Greece), Tyler Sorensen (UC Santa Cruz, USA), Simone Campanoni (Northwestern University, USA), Tor M. Aamodt (University of British Columbia, Canada), and David I. August (Princeton University, USA)*

*Yunzhe Liu (State Key Lab of Processors, Insititute of Computing Technology, CAS; University of Chinese Academy of Sciences), Xinyu Li (State Key Lab of Processors, Insititute of Computing Technology, CAS; University of Chinese Academy of Sciences), Tingting Zhang (Loongson Technology Co. Ltd.; Institute of Computing Technology, CAS), Tianyi Liu (The University of Texas at San Antonio), Qi Guo (State Key Lab of Processors, Insititute of Computing Technology, CAS), Fuxin Zhang (State Key Lab of Processors, Insititute of Computing Technology, CAS; University of Chinese Academy of Sciences), and Jian Wang (State Key Lab of Processors, Insititute of Computing Technology, CAS; University of Chinese Academy of Sciences)*

*Anubhav Bhatla (Indian Institute of Technology Bombay, India), Navneet Navneet (Indian Institute of Technology Bombay, India), and Biswabandan Panda (Indian Institute of Technology Bombay, India)*

# Session 1B: Emerging Technologies

# Session 2: Best Paper Session

## Session 3A: Networking

## Session 3B: Quantum Computing

## Session 4A: PIM Accelerators

## Session 4B: Cloud Technologies

## Session 5A: Tools and Analysis

## Session 5B: Accelerators for Emerging Workloads I

## Session 6A: NDP Technologies

## Session 6B: Security

## Session 6C: Parallel Architectures

## Session 7: Industry Session

## Session 8A: Machine Learning Accelerators I

## Session 8B: Compliers and Programming Models

## Session 9A: Machine Learning Accelerators II

## Session 9B: Memory Systems

# Session 10A: Prefetching

# Session 10B: Accelerators for Emerging Workloads II

**Author Index**