# 2021 ACM/IEEE 48th Annual International Symposium on Computer Architecture (ISCA 2021)

**Virtual Event**
**14 – 19 June 2021**

**Pages 1-566**

**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

# 2021 ACM/IEEE 48th Annual International Symposium on Computer Architecture (ISCA)

# ISCA 2021

## Table of Contents

## Session 1: Industry track

*Norman P. Jouppi (Google LLC), Doe Hyun Yoon (Google LLC), Matthew Ashcraft (Google LLC), Mark Gottscho (Google LLC), Thomas B. Jablin (Google LLC), George Kurian (Google LLC), James Laudon (Google LLC), Sheng Li (Google LLC), Peter Ma (Google LLC), Xiaoyu Ma (Google LLC), Thomas Norrie (Google LLC), Nishant Patil (Google LLC), Sushma Prasad (Google LLC), Clifford Young (Google LLC), Zongwei Zhou (Google LLC), and David Patterson (Google LLC)*

Brian W. Thompto (International Business Machines Corporation, Armonk,
USA), Dung Q. Nguyen (International Business Machines Corporation,
Armonk, USA), José E. Moreira (International Business Machines
Corporation, Armonk, USA), Ramon Bertran (International Business
Machines Corporation, Armonk, USA), Hans Jacobson (International
Business Machines Corporation, Armonk, USA), Richard J. Eickemeyer
(International Business Machines Corporation, Armonk, USA), Rahul M.
Rao (International Business Machines Corporation, Armonk, USA),
Michael Goulet (International Business Machines Corporation, Armonk,
USA), Marcy Byers (International Business Machines Corporation,
Armonk, USA), Christopher J. Gonzalez (International Business Machines
Corporation, Armonk, USA), Karthik Swaminathan (International Business
Machines Corporation, Armonk, USA), Nagu R. Dhanwada (International
Business Machines Corporation, Armonk, USA), Silvia M. Müller
(International Business Machines Corporation, Armonk, USA), Andreas
Wagner (International Business Machines Corporation, Armonk, USA),
Satish Kumar Sadasivam (International Business Machines Corporation,
Armonk, USA), Robert K. Montoye (International Business Machines
Corporation, Armonk, USA), William J. Starke (International Business
Machines Corporation, Armonk, USA), Christian G. Zoellin
(International Business Machines Corporation, Armonk, USA), Michael S.
Floyd (International Business Machines Corporation, Armonk, USA),
Jeffrey Stuecheli (International Business Machines Corporation,
Armonk, USA), Nandhini Chandramoorthy (International Business Machines
Corporation, Armonk, USA), John-David Wellman (International Business
Machines Corporation, Armonk, USA), Alper Buyuktosunoglu
(International Business Machines Corporation, Armonk, USA), Matthias
Pflanz (International Business Machines Corporation, Armonk, USA),
Balaram Sinharoy (International Business Machines Corporation, Armonk,
USA), and Pradip Bose (International Business Machines Corporation,
Armonk, USA)

Sukhan Lee (Memory Business Division, Samsung Electronics), Shin-haeng
Kang (Memory Business Division, Samsung Electronics), Jaehoon Lee
(Memory Business Division, Samsung Electronics), Hyeonsu Kim (Samsung
Advanced Institute of Technology, Samsung Electronics), Eojin Lee
(Memory Business Division, Samsung Electronics), Seungwoo Seo (Samsung
Advanced Institute of Technology, Samsung Electronics), Hosang Yoon
(Samsung Advanced Institute of Technology, Samsung Electronics),
Seungwon Lee (Samsung Advanced Institute of Technology, Samsung
Electronics), Kyounghwan Lim (Memory Business Division, Samsung
Electronics), Hyunsung Shin (Memory Business Division, Samsung
Electronics), Jinhyun Kim (Memory Business Division, Samsung
Electronics), Seongil O (Memory Business Division, Samsung
Electronics), Anand Iyer (Device Solutions America, Samsung
Electronics), David Wang (Device Solutions America, Samsung
Electronics), Kyomin Sohn (Memory Business Division, Samsung
Electronics), and Nam Sung Kim (Memory Business Division, Samsung
Electronics)

> Samuel Naffziger (Advanced Micro Devices, Inc.), Noah Beck (Advanced
> Micro Devices, Inc.), Thomas Burd (Advanced Micro Devices, Inc.),
> Kevin Lepak (Advanced Micro Devices, Inc.), Gabriel H. Loh (Advanced
> Micro Devices, Inc.), Mahesh Subramony (Advanced Micro Devices, Inc.),
> and Sean White (Advanced Micro Devices, Inc.)

## Session 2A: Microarchitecture-1

> Mainak Chaudhuri (Indian Institute of Technology Kanpur)

> Georgios Vavouliotis (Barcelona Supercomputing Center; Universitat
> Politècnica de Catalunya), Lluc Alvarez (Barcelona Supercomputing
> Center; Universitat Politècnica de Catalunya), Vasileios Karakostas
> (National Technical University of Athens), Konstantinos Nikas
> (National Technical University of Athens), Nectarios Koziris (National
> Technical University of Athens), Daniel A. Jiménez (Texas A&M
> University), and Marc Casas (Barcelona Supercomputing Center;
> Universitat Politècnica de Catalunya)

> Alberto Ros (University of Murcia, Spain ) and Alexandra Jimborean
> (University of Murcia, Spain)

## Session 2B: Memory-1

> Yifan Yuan (UIUC), Mohammad Alian (University of Kansas), Yipeng Wang
> (Intel), Ren Wang (Intel), Ilia Kurakin (Intel), Charlie Tai (Intel),
> and Nam Sung Kim (UIUC)

> Nezam Rohbani (Institute for Research in Fundamental Sciences (IPM),
> Iran), Sina Darabi (Sharif University of Technology, Iran), and Hamid
> Sarbazi-Azad (Institute for Research in Fundamental Sciences (IPM),
> Iran; Sharif University of Technology)

> Harini Muthukrishnan (University of Michigan), David Nellans (NVIDIA),
> Daniel Lustig (NVIDIA), Jeffrey A. Fessler (University of Michigan),
> and Thomas F. Wenisch (University of Michigan)

# Session 3A: Machine Learning-1

*Swagath Venkataramani (IBM Research, Yorktown Heights), Vijayalakshmi
Srinivasan (IBM Research, Yorktown Heights), Wei Wang (IBM Research,
Yorktown Heights), Sanchari Sen (IBM Research, Yorktown Heights),
Jintao Zhang (IBM Research, Yorktown Heights), Ankur Agrawal (IBM
Research, Yorktown Heights), Monodeep Kar (IBM Research, Yorktown
Heights), Shubham Jain (IBM Research, Yorktown Heights), Alberto
Mannari (IBM Research, Switzerland), Hoang Tran (IBM Research,
Yorktown Heights), Yulong Li (IBM Research, Yorktown Heights), Eri
Ogawa (IBM Research, Japan), Kazuaki Ishizaki (IBM Research, Japan),
Hiroshi Inoue (IBM Research, Japan), Marcel Schaal (IBM Research,
Yorktown Heights), Mauricio Serrano (IBM Research, Yorktown Heights),
Jungwook Choi (IBM Research, Yorktown Heights), Xiao Sun (IBM
Research, Yorktown Heights), Naigang Wang (IBM Research, Yorktown
Heights), Chia-Yu Chen (IBM Research, Yorktown Heights), Allison
Allain (IBM Research, Yorktown Heights), James Bonanno (IBM, Austin),
Nianzheng Cao (IBM Research, Yorktown Heights), Robert Casatuta (IBM,
Hopewell Junction), Matthew Cohen (IBM Research, Yorktown Heights),
Bruce Fleischer (IBM Research, Yorktown Heights), Michael Guillorn
(IBM Research, Yorktown Heights), Howard Haynie (IBM, Poughkeepsie),
Jinwook Jung (IBM Research, Yorktown Heights), Mingu Kang (IBM
Research, Yorktown Heights), Kyu-hyoun Kim (IBM Research, Yorktown
Heights), Siyu Koswatta (IBM Research, Yorktown Heights), Saekyu Lee
(IBM Research, Yorktown Heights), Martin Lutz (IBM Research, Yorktown
Heights), Silvia Mueller (IBM, Germany), Jinwook Oh (IBM Research,
Yorktown Heights), Ashish Ranjan (IBM Research, Yorktown Heights),
Zhibin Ren (IBM Research, Yorktown Heights), Scot Rider (IBM,
Poughkeepsie), Kerstin Schelm (IBM, Germany), Micheal Scheuermann (IBM
Research, Yorktown Heights), Joel Silberman (IBM Research, Yorktown
Heights), Jie Yang (IBM Research, Yorktown Heights), Vidhi Zalani (IBM
Research, Yorktown Heights), Xin Zhang (IBM Research, Yorktown
Heights), Ching Zhou (IBM Research, Yorktown Heights), Matt Ziegler
(IBM Research, Yorktown Heights), Vinay Shah (IBM, United Kingdom),
Moriyoshi Ohara (IBM Research, Japan), Pong-Fei Lu (IBM Research,
Yorktown Heights), Brian Curran (IBM, Poughkeepsie), Sunil Shukla (IBM
Research, Yorktown Heights), Leland Chang (IBM Research, Yorktown
Heights), and Kailash Gopalakrishnan (IBM Research, Yorktown Heights)*

*Anant Nori (Processor Architecture Research Lab, Intel Labs, Intel),
Rahul Bera (ETH Zurich), Shankar Balachandran (Processor Architecture
Research Labs Intel), Joydeep Rakshit (Processor Architecture Research
Labs Intel), Omer Om (Processor Architecture Research Labs Intel),
Avishaii Abuhatzera (Intel), Belliappa Kuttanna (Intel), and Sreenivas
Subramoney (Processor Architecture Research Labs Intel)*

*Jiayi Huang (UC Santa Barbara), Pritam Majumder (Texas A&M
University), Sungkeun Kim (Texas A&M University), Abdullah Muzahid
(Texas A&M University), Ki Hwan Yum (Texas A&M University), and Eun
Jung Kim (Texas A&M University)*

## Session 3B: Microarchitecture-2

## Session 4A: Processing in/near Memory

## Session 4B: Data Center

## Session 5A: Security-1

## Session 5B: Accelerators-1

## Session 6A: Compilers

## Session 6B: Memory-2

## Session 7A: Accelerators-2

## Session 7B: Graph Processing

## Session 8A: Low Temperature/Low Energy Computing

## Session 8B: Machine Learning-2

## Session 9A: Memory-3

## Session 9B: Network Storage and Acceleration

## Session 10A: Quantum/Photonics

## Session 10B: Reliability and Security

## Session 11A: DRAM/IO/Network

## Session 11B: Security-2

## Session 12A: Accelerators-3

## Session 12B: Sparse Processing

**Author Index**