# 2022 IEEE International Symposium on High-Performance Computer Architecture (HPCA 2022)

**Virtual Symposium**
**2-6 April 2022**

**Pages 1-618**

**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

# 2022 IEEE International Symposium on High-Performance Computer Architecture (HPCA)
# HPCA 2022

## Table of Contents

## Session 1A: Accelerators I

## Session 1B: Security I

## Session 1C: At Scale

## Session 2A: Accelerators II

## Session 2B: Security II

# Session 2C: Quantum I

# Session 3A: Accelerators III

# Session 3B: Security III

## Session 3C: Quantum II

## Session 4A: Accelerators IV

## Session 4B: Storage, Scheduling, Interfaces

## Session 4C: Best Paper Candidates

## Session 5A: Simulation

## Session 5B: Cache Hierarchy

## Session 5C: Quantum III

## Session 6A: Synthesis

## Session 6B: Traditional Architecture

## Session 7A: Accelerators V

## Session 7B: Non-Volatile Memory

## Session 7C: Network On Chip

# Session 8A: Accelerators VI

# Session 8B: Memory

## Session 8C: Industrial Session

**Author Index**