# 2021 30th International Conference on Parallel Architectures and Compilation Techniques (PACT 2021)

Atlanta, Georgia, USA
26-29 September 2021

**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

# 2021 30th International Conference on Parallel Architectures and Compilation Techniques (PACT)
# PACT 2021

## Table of Contents

## Session 1: Tuning and Lifting

*Phitchaya Mangpo Phothilimthana (Google, USA), Amit Sabne (Google, USA), Nikhil Sarda (Google, USA), Karthik Srinivasa Murthy (Google, USA), Yanqi Zhou (Google, USA), Christof Angermueller (Google, USA), Mike Burrows (Google, USA), Sudip Roy (Google, USA), Ketan Mandke (Google, USA), Rezsa Farahani (Google, USA), Yu Emma Wang (Google, USA), Berkin Ilbeyi (Google, USA), Blake Hechtman (Google, USA), Bjarke Roune (Google, USA), Shen Wang (Google, USA), Yuanzhong Xu (Google, USA), and Samuel J. Kaufman (University of Washington, USA)*

*Alexander Brauckmann (TU Dresden, Germany), Andrés Goens (Barkhausen Institut, Dresden, Germany), and Jeronimo Castrillon (TU Dresden, Germany)*

*Geonhwa Jeong (Georgia Institute of Technology), Gokcen Kestor (Pacific Northwest National Laboratory), Prasanth Chatarasi (IBM Research), Angshuman Parashar (NVIDIA), Po-An Tsai (NVIDIA), Sivasankaran Rajamanickam (Sandia National Laboratories), Roberto Gioiosa (Pacific Northwest National Laboratory), and Tushar Krishna (Georgia Institute of Technology)*

*William Moses (MIT CSAIL, USA), Lorenzo Chelini (TU Eindhoven, The Netherlands), Ruizhe Zhao (Imperial College London, UK), and Oleksandr Zinenko (Google Inc., France)*

*Bruce Collie (University of Edinburgh, United Kingdom) and Michael O'Boyle (University of Edinburgh, United Kingdom)*

## Session 2: Heterogeneous Systems

*Joonsung Kim (Seoul National University, Republic of Korea), Suyeon
Hur (Seoul National University, Republic of Korea), Eunbok Lee (Seoul
National University, Republic of Korea), Seungho Lee (Seoul National
University, Republic of Korea), and Jangwoo Kim (Seoul National
University, Republic of Korea)*

*Myeonggyun Han (UNIST) and Woongki Baek (UNIST)*

*Naveen Vedula (Simon Fraser University, Canada), Reza Hojabr (Simon
Fraser University, Canada; University of Tehran, Iran), Ahmad Khonsari
(University of Tehran, Iran; IPM, Tehran, Iran), and Arrvindh
Shriraman (Simon Fraser University, Canada)*

*Daehyeon Baek (KAIST, South Korea), Soojin Hwang (KAIST, South Korea),
Taekyung Heo (KAIST, South Korea), Daehoon Kim (DGIST, South Korea),
and Jaehyuk Huh (KAIST, South Korea)*

*Maximilian Lam (Harvard University, USA), Zachary Yedidia (Harvard
University, USA), Colby R. Banbury (Harvard University, USA), and
Vijay Janapa Reddi (Harvard University, USA)*

## Session 3: Characterization and Near-Memory Computing

*Wanling Gao (Institute of Computing Technology, Chinese Academy of
Sciences; BenchCouncil (International Open Benchmark Council);
University of Chinese Academy of Sciences), Fei Tang (Institute of
Computing Technology, Chinese Academy of Sciences; University of
Chinese Academy of Sciences), Jianfeng Zhan (Institute of Computing
Technology, Chinese Academy of Sciences; BenchCouncil (International
Open Benchmark Council); University of Chinese Academy of Sciences),
Xu Wen (Institute of Computing Technology, Chinese Academy of
Sciences; University of Chinese Academy of Sciences), Lei Wang
(Institute of Computing Technology, Chinese Academy of Sciences;
BenchCouncil (International Open Benchmark Council); University of
Chinese Academy of Sciences), Zheng Cao (Alibaba), Chuanxin Lan
(Institute of Computing Technology, Chinese Academy of Sciences),
Chunjie Luo (Institute of Computing Technology, Chinese Academy of
Sciences; BenchCouncil (International Open Benchmark Council);
University of Chinese Academy of Sciences), Xiaoli Liu (Alibaba), and
Zihan Jiang (Institute of Computing Technology, Chinese Academy of
Sciences; University of Chinese Academy of Sciences)*

## Session 4: Memory Hierarchy

## Session 5: Graphs and Applications