# 2021 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS 2021)

**Virtual Conference**
**28 – 30 March 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

CURRAN ASSOCIATES INC.
proceedings
.com

# 2021 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)

# ISPASS 2021

## Table of Contents

## Paper Session I: Benchmarking

Arun Subramaniyan (University of Michigan, USA), Yufeng Gu (University
of Michigan, USA), Timothy Dunn (University of Michigan, USA), Somnath
Paul (Intel Corporation, USA), Md Vasimuddin (Intel Corporation,
India), Sanchit Misra (Intel Corporation, India), David Blaauw
(University of Michigan, USA), Satish Narayanasamy (University of
Michigan, USA), and Reetuparna Das (University of Michigan, USA)

Trinayan Baruah (Northeastern University), Kaustubh Shivdikar
(Northeastern University), Shi Dong (Cerebras Systems), Yifan Sun
(William & Mary), Saiful A Mojumder (Boston University), Kihoon Jung
(KAIST), José L. Abellán (Universidad Católica de Murcia), Yash
Ukidave (AMD), Ajay Joshi (Boston University), John Kim (KAIST), and
David Kaeli (Northeastern University)

# Paper Session II: GPUs

## Poster Session A:

## Paper Session III: Characterization

## Paper Session IV: Software Analysis

## Paper Session V: Best Paper Nominations

## Poster Session B:

## Paper Session VI: Datacenters and HPC

## Paper Session VII: HW and Co-Design