# 2025 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS 2025)

**Ghent, Belgium**
**11-13 May 2025**

**Additional Copies of This Publication Are Available From:**

CURRAN ASSOCIATES INC.
proceedings
.com

# 2025 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)

# ISPASS 2025

## Table of Contents

## Machine Learning

*Rachid Karami (University of California, USA), Sheng-Chun Kao
(Mountain View, USA), and Hyoukjun Kwon (University of California,
USA)*

*Junsoo Kim (HyperAccel Inc., South Korea), Hunjong Lee (HyperAccel
Inc., South Korea), Geonwoo Ko (KAIST, South Korea), Gyubin Choi
(HyperAccel Inc., South Korea), Seri Ham (KAIST, South Korea),
Seongmin Hong (HyperAccel Inc., South Korea), and Joo-Young Kim
(HyperAccel Inc., South Korea)*

*Zishen Wan (Georgia Institute of Technology), Jiayi Qian (Georgia
Institute of Technology), Yuhang Du (University of Minnesota, Twin
Cities), Jason Jabbour (Harvard University), Yilun Du (Harvard
University), Yang Katie Zhao (University of Minnesota, Twin Cities),
Arijit Raychowdhury (Georgia Institute of Technology), Tushar Krishna
(Georgia Institute of Technology), and Vijay Janapa Reddi (Harvard
University)*

*Jamin Seo (Georgia Institute of Technology, USA), Jianming Tong
(Georgia Institute of Technology, USA), Tushar Krishna (Georgia
Institute of Technology, USA), and Hyoukjun Kwon (University of
California, USA)*

## GPU/Ray Tracing

## System Characterization and Framework

## Benchmarking and Tools

## Modeling and Evaluation

## Accelerators and Power

## Posters