# 2024 IEEE 6th International Conference on AI Circuits and Systems (AICAS 2024)

**22-25 April 2024**

**Abu Dhabi, United Arab Emirates**

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

# AICAS 2024 TABLE OF CONTENTS

## A1L-A
## AI and ML at the Edge

Date:    Tuesday, April 23, 2024
Time:    9:00 - 10:30
Room:    B 00055
Chair:    Gianluca Setti, KAUST

Date: Tuesday, April 23, 2024
Time: 9:00 - 10:30
Room: B 00107
Chair: Zihang Song, King's College London

# A1L-C
# ML Applications in Communications and ASR

Date:     Tuesday, April 23, 2024
Time:     9:00 - 10:30
Room:     KIC
Chair:    Alex James, Digital University Kerala

# A2L-A
# Hardware Accelerator ICs

Date: Tuesday, April 23, 2024
Time: 13:10 - 14:40
Room: B 00055
Chair: Mahmoud Al-Qutayri, Khalifa University

# A2L-B
# Spiking Neural Networks II

Date:     Tuesday, April 23, 2024
Time:     13:10 - 14:40
Room:     B 00107
Chair:    Hani Saleh, Khalifa University (TBC)

# A2L-C
# Security, Safety and Trustworthiness

Date:     Tuesday, April 23, 2024
Time:     13:10 - 14:40
Room:     KIC
Chair:    Abdulhadi Shoufan, Khalifa University

# A3L-A
# Processors for AI

Date: Tuesday, April 23, 2024
Time: 15:00 - 16:30
Room: B 00055
Chair: Thanos Stouraitis, Khalifa University

# A3L-B
# In-memory Computing

Date: Tuesday, April 23, 2024
Time: 15:00 - 16:30
Room: B 00107
Chair: Khaled Nabil Salama, KAUST (TBC)

# A3L-C
# Object Detection and Applications

Date:      Tuesday, April 23, 2024
Time:      15:00 - 16:30
Room:      KIC
Chair:     Naoufel Werghi, Khalifa University

# B1L-A
# FPGA-based Accelerators

Date:    Wednesday, April 24, 2024
Time:    9:00 - 10:30
Room:   B 00055
Chair:   Abdulhadi Shoufan, Khalifa University

# B1L-B
# Analog Accelerators

Date:    Wednesday, April 24, 2024
Time:    9:00 - 10:30
Room:    B 00107
Chair:   Mihai Sanduleanu, Khalifa University

# B1L-C
# Anomaly Detection Techniques and Applications

Date: Wednesday, April 24, 2024
Time: 9:00 - 10:30
Room: KIC
Chair: Khaled Nabil Salama, KAUST

# B2L-A
# Flexible ML Architectures

Date:     Wednesday, April 24, 2024
Time:     10:50 - 12:20
Room:     B 00055
Chair:    Terry Tao Ye, Southern University of Science and Technology

# B2L-B
# Neuromorphic Systems

Date: Wednesday, April 24, 2024
Time: 10:50 - 12:20
Room: B 00107
Chair: Mohamad Sawan, Westlake University

# B2L-C
# Computer Vision Algorithms and Architectures

Date:     Wednesday, April 24, 2024
Time:     10:50 - 12:20
Room:     KIC
Chair:    Naoufel Werghi, Khalifa University

# B3L-A
# Arithmetic for Machine Learning

Date:     Wednesday, April 24, 2024
Time:     14:20 - 15:50
Room:    B 00055
Chair:    Shervin Vakili, INRS-EMT University

Ming-Guang Lin, Jiing-Ping Wang, Cheng-Yang Chang, An-Yeu Wu

*National Taiwan University, Taiwan*

Shervin Vakili

*INRS-EMT University, Canada*

Zhibin Luo, Junyi Mai, Enyi Yao

*South China University of Technology, China*

Junnosuke Suzuki, Mari Yasunaga, Kazushi Kawamura, Thiem Van Chu, Masato Motomura

*Tokyo Institute of Technology, Japan*

Rishi Agrawal[1], Narayanabhatla Savyasachi Abhijith[1], Uppugunduru Anil Kumar[2], Sreehari Veeramachaneni[3], Syed Ershad Ahmed[1]

*[1]Birla Institute of Technology and Science, Pilani, India; [2]ICFAI Foundation for Higher Education, India; [3]Gokaraju Rangaraju Institute of Engineering and Technology, India*

# B3L-B
# Memristors for ML/AI

Date:     Wednesday, April 24, 2024
Time:     14:20 - 15:50
Room:   B 00107
Chair:    Bhaskar Choubey, Siegen University

Brady Taylor[1], Xiaoxuan Yang[1,2], Hai Li[1]

[1]*Duke University, United States; *[2]*University of Virginia, United States*

Sumit Diware[1], Mohammad Amin Yaldagard[1], Anteneh Gebregiorgis[1], Rajiv V. Joshi[2], Said Hamdioui[1], Rajendra Bishnoi[1]

[1]*Delft University of Technology, Netherlands; *[2]*IBM Thomas J. Watson Research Centre, United States*

Zhenming Yu[1,2], Ming-Jay Yang[1,2], Jan Finkbeiner[1,2], Sebastian Siegel[2], John Paul Strachan[1,2], Emre Neftci[1,2]

[1]*RWTH Aachen, Germany; *[2]*Forschungszentrum Jülich GmbH, Germany*

Pierre Lewden[1], Adrien F. Vincent[2], Jean Tomas[2], Chip-Hong Chang[3], Sylvain Saïghi[1,2]

[1]*CNRS@CREATE, Singapore; *[2]*Université de Bordeaux, France; *[3]*Nanyang Technological University, Singapore*

Zidu Li, Phil David Börner, Maurice Müller, Andreas Bablich, Peter Haring Bolívar, Bhaskar Choubey

*Universität Siegen, Germany*

# B3L-C
# Convolutional Neural Networks

Date:      Wednesday, April 24, 2024
Time:      14:20 - 15:50
Room:      KIC
Chair:     Vasilis Sakellariou, Khalifa University

# B4L-A
# Quantization in Machine Learning Systems

Date: Wednesday, April 24, 2024
Time: 16:10 - 17:40
Room: B 00055
Chair: Vassilis Paliouras, University of Patras

# B4L-B
# Implementation Techniques

Date:     Wednesday, April 24, 2024
Time:     16:10 - 17:40
Room:     B 00107
Chair:    José M. de la Rosa, Institute of Microelectronics of Seville

Shreyas Deshmukh, Shubham Patil, Anmol Biswas, Vivek Saraswat, Abhishek Kadam, Ajay K. Singh,
Laxmeesha Somappa, Maryam Shojaei Baghini, Udayan Ganguly

*Indian Institute of Technology Bombay, India*

I-Chun Liu[1], Chun-Jui Chen[1], Xiu-Zhu Li[1], Yong-Qi Cheng[1], Chung-Wei Huang[1], Pin-Han Lin[1],
Hsuan-Wei Pu[1], Sheng-Yu Peng[1], Yu Tsao[2]

*[1]National Taiwan University of Science and Technology, Taiwan; [2]Academia Sinica, Taiwan*

Gilha Lee, Seungil Lee, Hyun Kim

*Seoul National University of Science and Technology, Korea*

Surajit Bhattacherjee, Daksh Shah, Dipankar Pal

*Birla Institute of Technology and Science, Pilani, India*

Bhaskar Choubey, Hendrik Sommerhoff, Michael Moeller, Andreas Kolb

*Universität Siegen, Germany*

# B4L-C
# Large Language Models

Date:     Wednesday, April 24, 2024
Time:     16:10 - 17:40
Room:     KIC
Chair:    Xuan-Truong Nguyen, Seoul National University

Janak Sharda, Po-Kai Hsu, Shimeng Yu

*Georgia Institute of Technology, United States*

Khaleelulla Khan Nazeer[1], Mark Schöne[1], Rishav Mukherji[2], Bernhard Vogginger[1], Christian Mayr[1], David Kappel[3], Anand Subramoney[4]

*[1]Technische Universität Dresden, Germany; [2]Birla Institute of Technology and Science, Pilani, India; [3]Ruhr Universität Bochum, Germany; [4]University of London, United Kingdom*

Minseok Seo, Seongho Jeong, Hyuk-Jae Lee, Xuan-Truong Nguyen

*Seoul National University, Korea*

You-En Wu[1], Hsin-I Wu[2], Kuo-Cheng Chin[2], Yi-Chun Yang[3], Ren-Song Tsay[3]

*[1]Taipei Fuhsing Private School, Taiwan; [2]DeepMentor, Taiwan; [3]National Tsing Hua University, Taiwan*

Salah Eddine Bekhouche[1], Abdenour Hadid[2]

*[1]University of the Basque Country UPV/EHU, Spain; [2]Sorbonne University Abu Dhabi, U.A.E.*

# C1L-A
# Special Session: Intersection of Hardware and Software for Edge AI and TinyML

Date:     Thursday, April 25, 2024
Time:     9:00 - 10:30
Room:     B 00055
Chairs:   Wei Mao, Xidian University
          Juntao Guan, Xidian University

## ILD-MPQ: Learning-Free Mixed-Precision Quantization with Inter-layer Dependency Awareness

Ruge Xu[1], Qiang Duan[2], Qibin Chen[2], Xinfei Guo[1]

[1]Shanghai Jiao Tong University, China; [2]Inspur Academy of Science and Technology, China

## Q-Segment: Segmenting Images In-Sensor for Vessel-Based Medical Diagnosis

Pietro Bonazzi, Yawei Li, Sizhen Bian, Michele Magno

ETH Zürich, Switzerland

## Microarchitecture Aware Neural Architecture Search for TinyML Devices

Juntao Guan, Gufeng Liu, Fanhong Zeng, Rui Lai, Ruixue Ding, Zhangming Zhu

Xidian University, China

## An Energy-Efficient Look-Up Table Framework for Super Resolution on FPGA

Huanan Li, Shicheng Jia, Juntao Guan, Rui Lai, Shubin Liu, Zhangming Zhu

Xidian University, China

## A Hardware-Efficient EMG Decoder with an Attractor-Based Neural Network for Next-Generation Hand Prostheses

Mohammad Kalbasi[1,2], MohammadAli Shaeri[1], Vincent Alexandre Mendez[1], Solaiman Shokur[1], Silvestro Micera[1], Mahsa Shoaran[1]

[1]École Polytechnique Fédérale de Lausanne, Switzerland; [2]Sharif University of Technology, Iran

# C1L-B
# Design Techniques for ML Hardware

Date:     Thursday, April 25, 2024
Time:     9:00 - 10:30
Room:     B 00107
Chair:    Hani Saleh, Khalifa University

# C1L-C
# Health and Brain-Machine Applications

# C2L-A
# Circuits and Systems for AI

Date:     Thursday, April 25, 2024
Time:     10:50 - 12:20
Room:     B 00055
Chair:    Tzi-Dar Chiueh, National Taiwan University

Pi-Chuan Chen[1], Yu-Tung Liu[2], Guo-Yang Zeng[2], Tzi-Dar Chiueh[1]

*[1]National Taiwan University, Taiwan; [2]Mediatek Inc., Taiwan*

Yuhan Qin, Yulong Meng, Haitao Du, Yazhuo Guo, Yi Kang

*University of Science and Technology of China, China*

Hongbo Guo[1,2,3], Hao Wu[2,3], Jiahui Xia[4], Yibang Cheng[4], Qianhui Guo[4], Yi Chen[4], Tingyan Xu[4], Jiguang Wang[4], Guoxing Wang[1,2,3]

*[1]Shanghai Jiao Tong University, China; [2]RingConn LLC, United States; [3]Shenzhen Ninenovo Technology Limited, China; [4]Shanghai Jiao Tong University School of Medicine, China*

A. Bazzi, E. Hardy, J. Ballester, F. Badets, L. Hutin

*CEA-Leti, Université Grenoble Alpes, France*

Paulo Machado[1,2], Ruxandra Barbulescu[1], Luis Miguel Silveira[1,2]

*[1]INESC-ID, Portugal; [2]Universidade de Lisboa, Portugal*

# Additional Paper:

Junfeng Tan, Guosheng Yu, Jianing Li, Xiaohan Ma, Fang Bao, Evens Pan, David Bian, Yongfu Li, Yuan Du, Li Du, Bo Li, Wei Mao