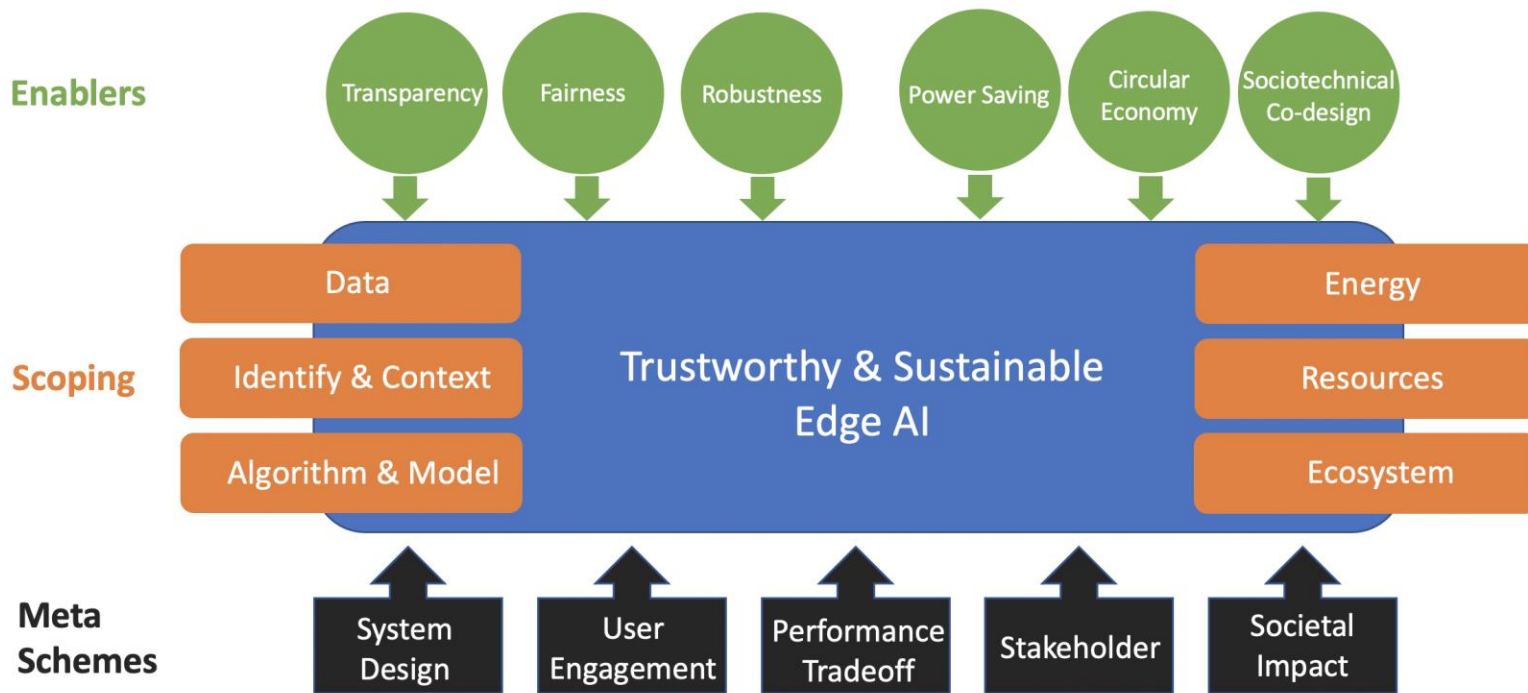# *Sustainable & Trustworthy*

# *Edge AI for Future Computing*

**Aaron Ding**

Director of CPI Lab, TU Delft
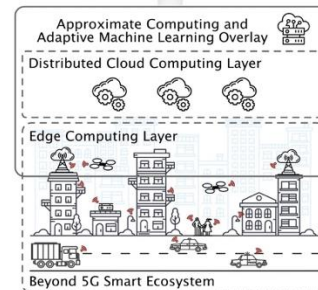
# Complex Subject



Aaron Ding, Marijn Janssen, Jon Crowcroft. "**Trustworthy and Sustainable Edge AI: A Research Agenda**"

# CPI on Edge AI



- SPATIAL of €5M grant
- APROPOS of €4M grant

**Trustworthy Edge AI**

**Score: 98/100 | Rate: 8%**



**Sustainable Edge AI**

**Score: 14.5/15 | Rate: 3%**



**TU**Delft

3

# Is **Edge AI** a Real Thing?

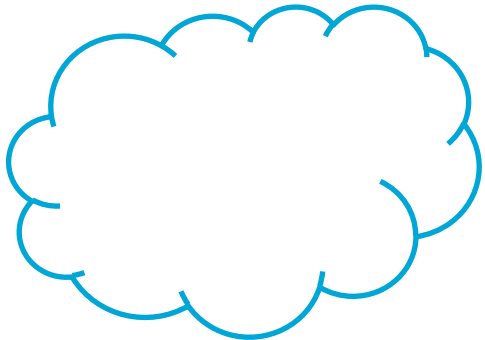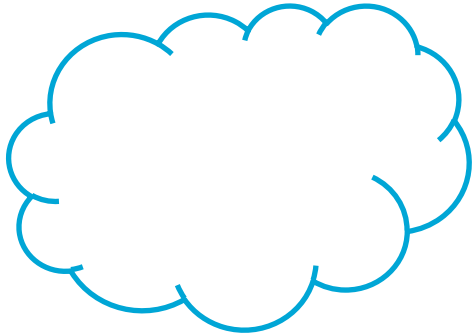# Edge AI is **Real**

| Provider | Hardware |
|----------|----------|
| Google | Tensor Processing Units (TPU) |
| Intel | Movidius Vision Processing Units (VPUs) & Xeon D-2100 |
| Qualcomm | Qualcomm Snapdragon 8 Series, Hexagon DSP |
| Huawei | Ascend Series & Kirin 600/900 NPU |
| Samsung | Exynos 9820 Neural Processing Unit (NPU) |
| NVIDIA | TURING GPU |

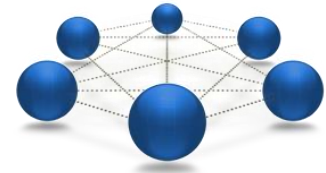| Provider | Dev Platform |
|----------|--------------|
| Microsoft | Azure Data Box Edge |
| Intel | Movidius Neural Compute Stick |
| NVIDIA | Jetson Nano, TX, Xavier NX |
| Huawei | Atlas AI Computing Platform |

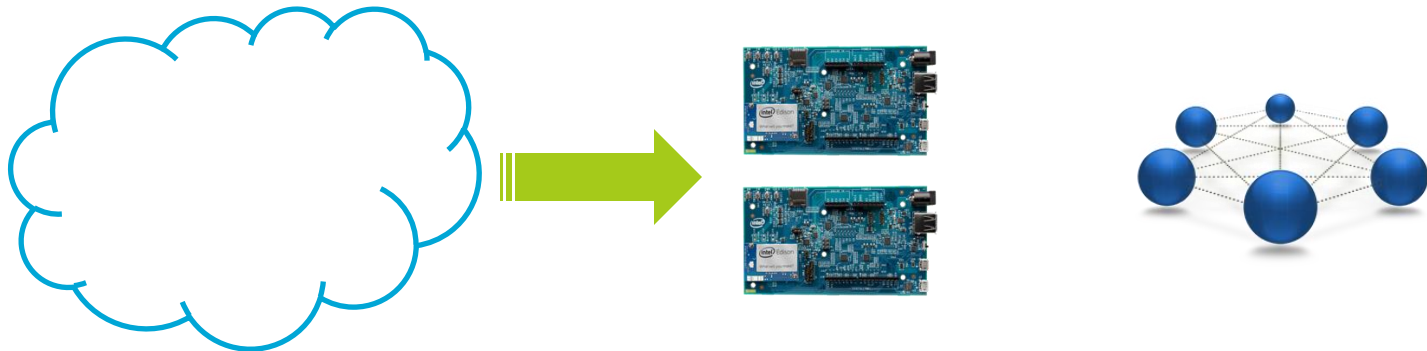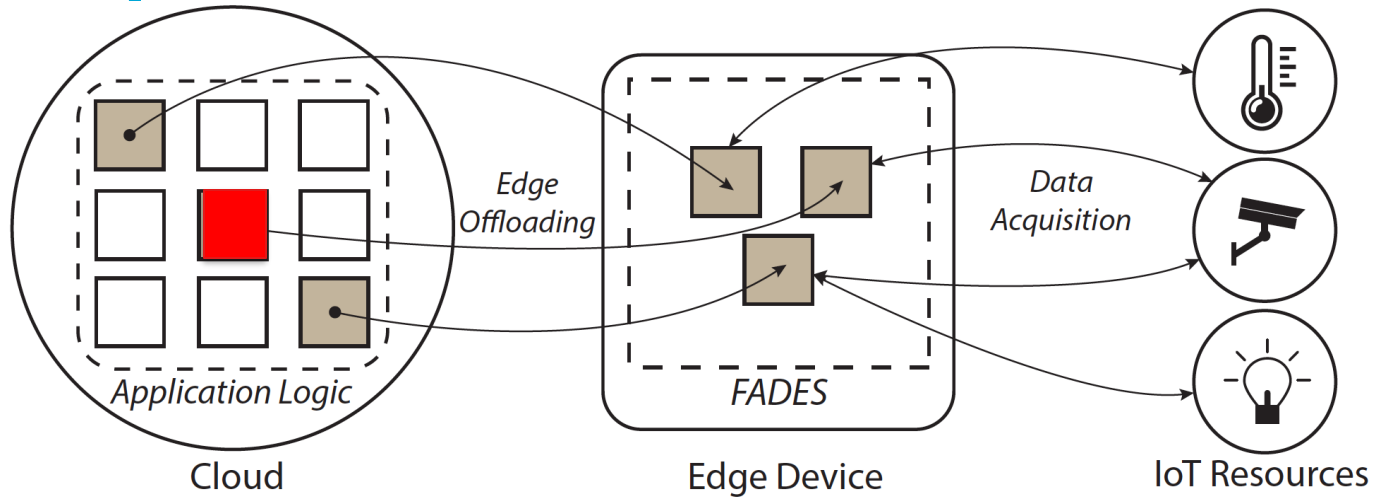| Provider | Management Framework |
|----------|----------------------|
| Microsoft | Azure IoT Edge |
| Google | Google Cloud IoT |
| NVIDIA | NVIDIA EGX |
| Amazon | AWS IoT Greengrass |
| Alibaba | Link IoT Edge |
| Linux Found. | EdgeX & Akraino Edge Stack |
| Huawei | KubeEdge |

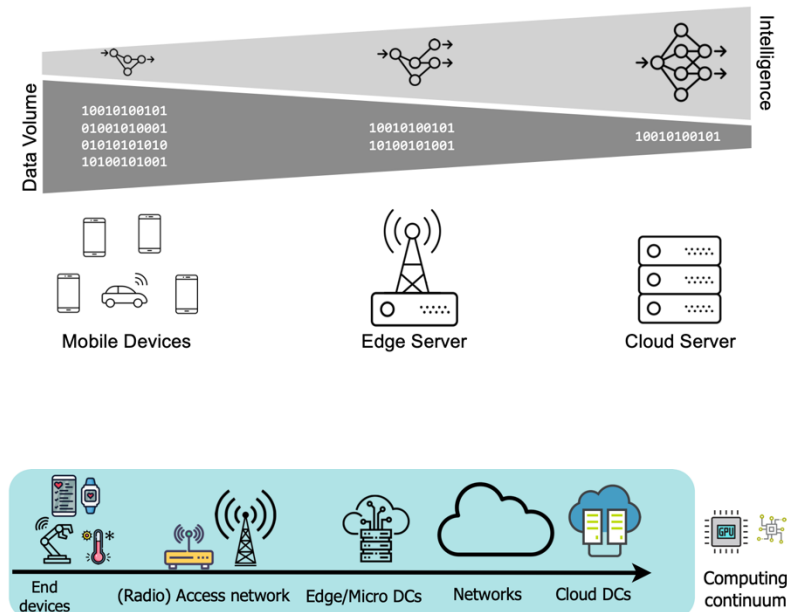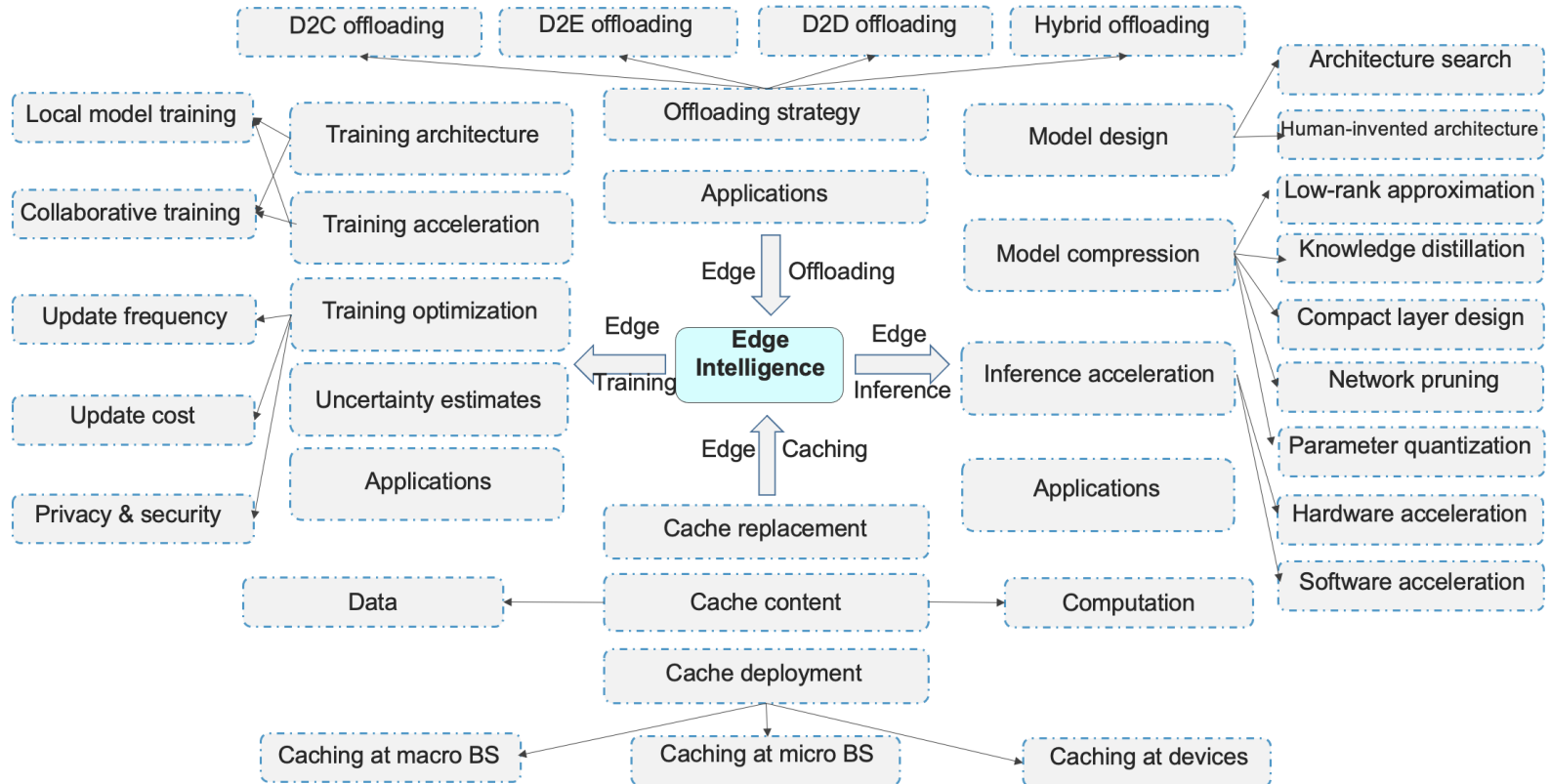# *What exactly is Edge AI ?*

# Edge Paradigm



Edge

IoT

# Example



Edge Offloading

Data Acquisition

Application Logic

Cloud

FADES

Edge Device

IoT Resources

# Continuum Perspective

**Degree of Decentralization** (↑)

**Potential Data Offloaded to Cloud** (↓)

Level 6: On-device Training & Inference

Level 5: Edge Training & Edge Inference

Level 4: Cloud-Edge Co-training & Co-Inference

Level 3: Cloud Training & On-device Inference

Level 2: Cloud Training & Edge Inference

Level 1: Cloud Training & Cloud-Edge Co-inference

Level 0: Fully Cloud Training & Inference



Intelligence

Data Volume

10010100101
01001010001
0101010101010
10100101001

10010100101
10100101001

10010100101

Mobile Devices

Edge Server

Cloud Server



End devices

(Radio) Access network

Edge/Micro DCs

Networks

Cloud DCs

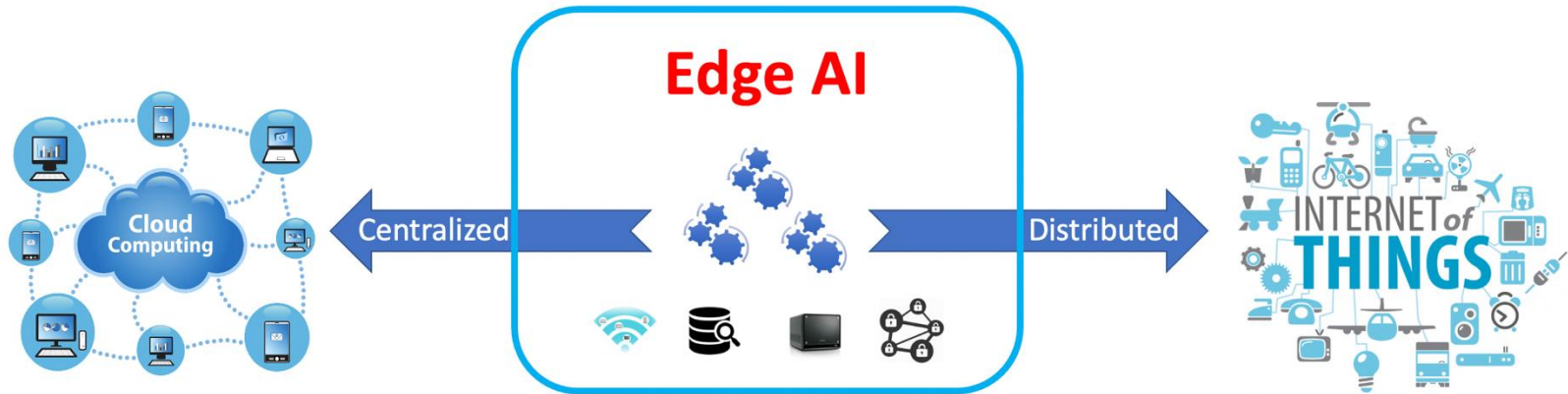Computing continuum

# Enabling Techniques

# Motivation

Bridge the Gap



## Consolidate Cloud & IoT

# Case: Crowd Intelligence on Edge

- **<u>Societal impact</u>** of past years
  - Responding and coping with emergency/pandemics
  - <span style="color:red">Urban activity/mobility sensing</span> on the edge

# **Motivation**

high fidelity entails high cost, infrastructure dependency, privacy intervention

- Low cost and scalable
  - User equipment
  - Deployment and coverage

- Passive (non-intrusive)
  - No need to force user interactions nor mandatory engagement

- Privacy-aware/friendly
  - Balance fidelity and data (local) regulations

# Unexpected

- Project ends…
  - Regulatory and legal considerations
  - Privacy in local context

# Lessons

- Boundary + Awareness
  - Privacy on Edge?  regulations and legal
  - Difference across countries

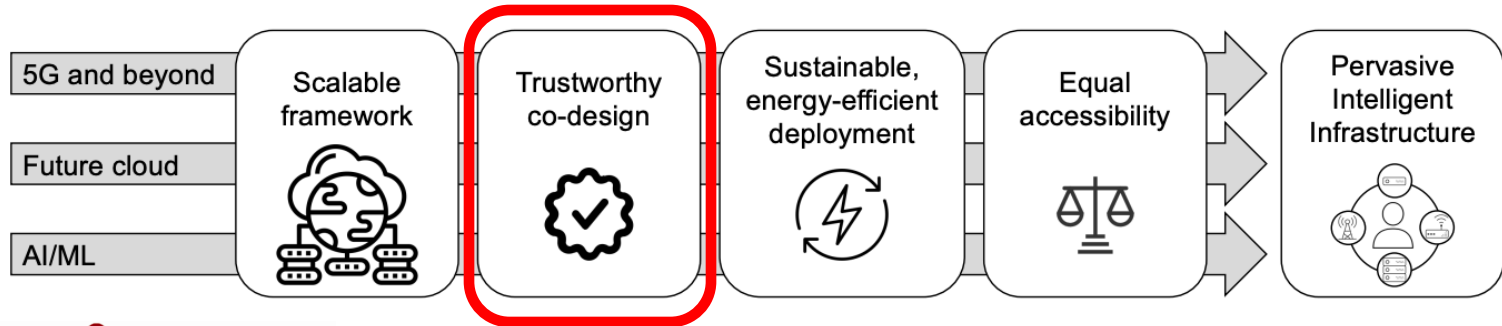## Where is my Tag? Unveiling Alternative Uses of Apple FindMy Service

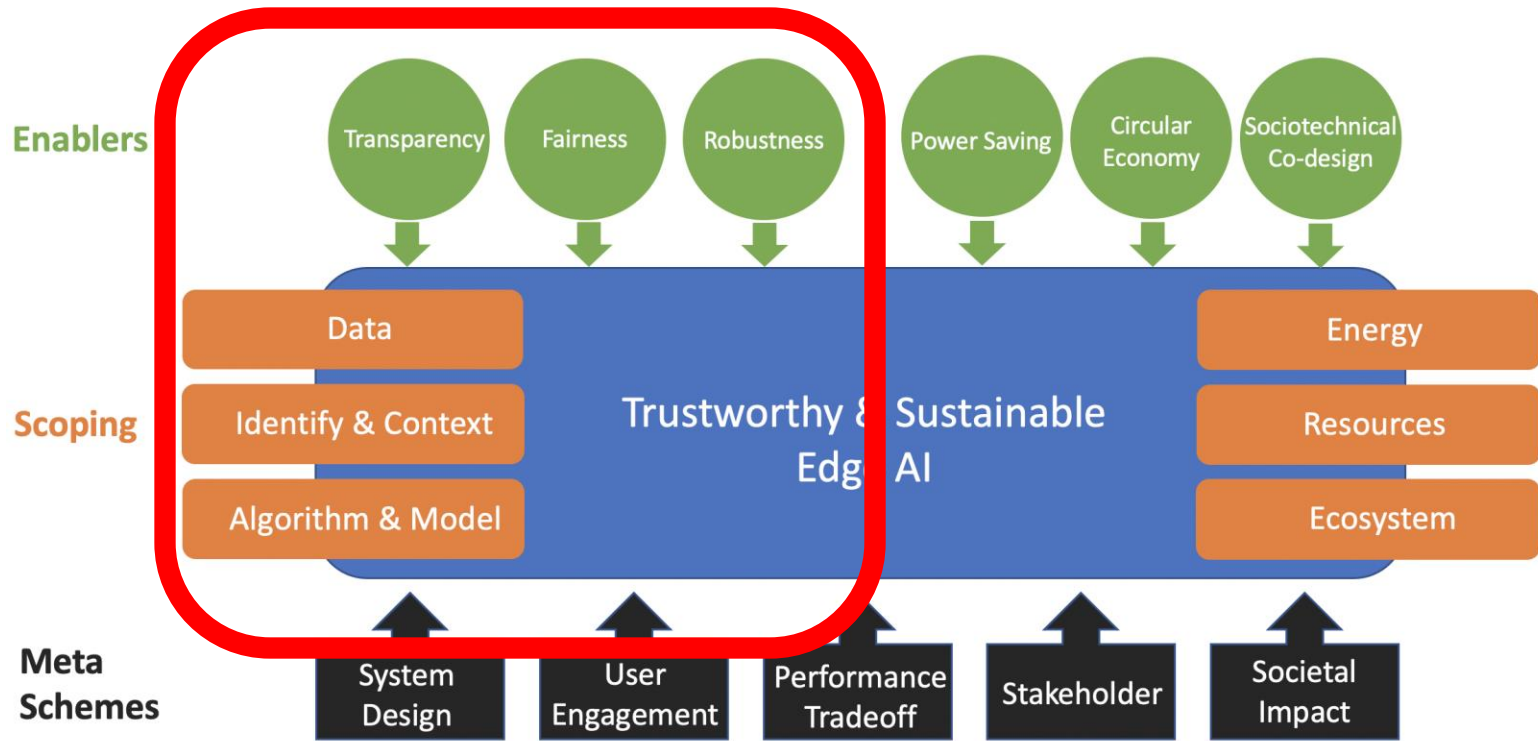"Learn from the mistakes of others. You can't live long enough to make them all yourself." - Eleanor Roosevelt

**TU**Delft

# Roadmap



**Roadmap for Edge AI: A Dagstuhl Perspective**

Aaron Yi Ding[1]*, Ella Peltonen[2], Tobias Meuser[3], Atakan Aral[4], Christian Becker[5], Schahram Dustdar[6], Thomas Hiessl[6], Dieter Kranzlmüller[7], Madhusanka Liyanage[8], Setareh Magshudi[9], Nitinder Mohan[10], Jörg Ott[10], Jan S. Rellermeyer[11,1], Stefan Schulte[12], Henning Schulzrinne[13], Gürkan Solmaz[14], Sasu Tarkoma[15], Blesson Varghese[16], Lars Wolf[17]

ACM SIGCOMM CCR, Vol. 52, No.1, 2022
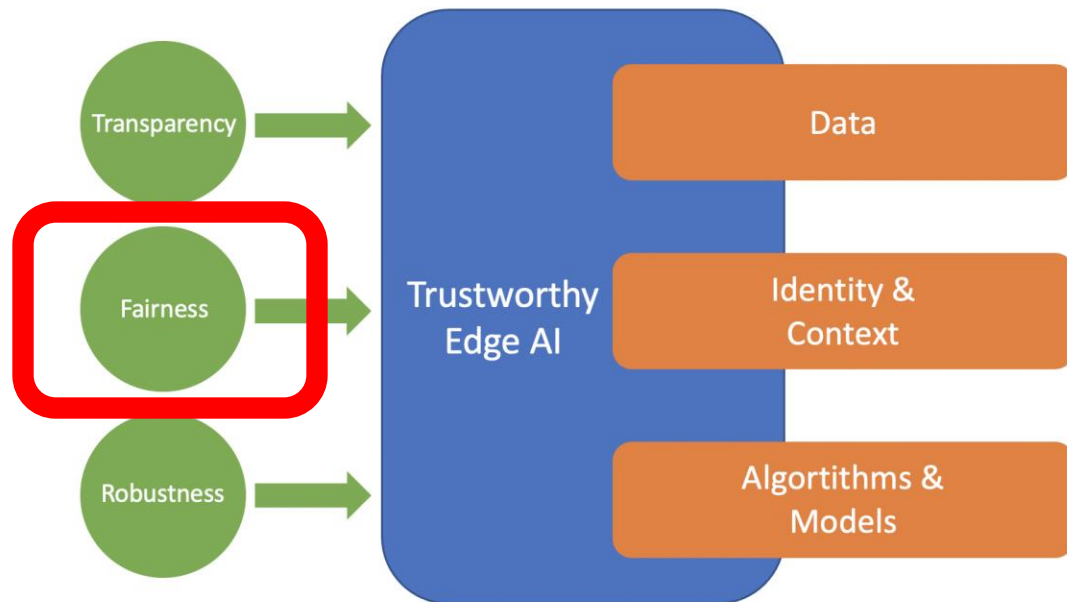
# Research Agenda



Aaron Ding, Marijn Janssen, Jon Crowcroft. "Trustworthy and Sustainable Edge AI: A Research Agenda", IEEE TPS

# Targets

- Enabler 🟢
  - Transparency
  - Fairness
  - Robustness

- Scope 🟠
  - Data
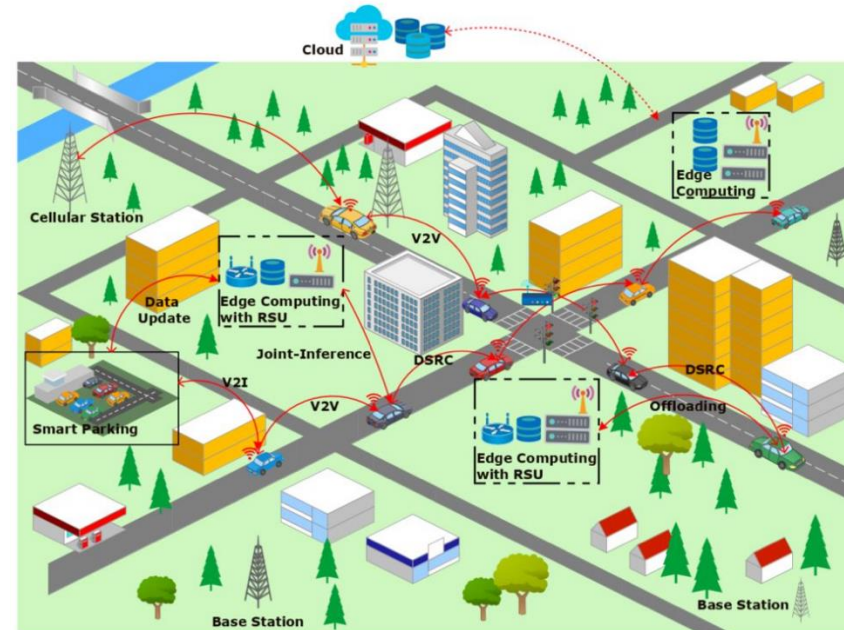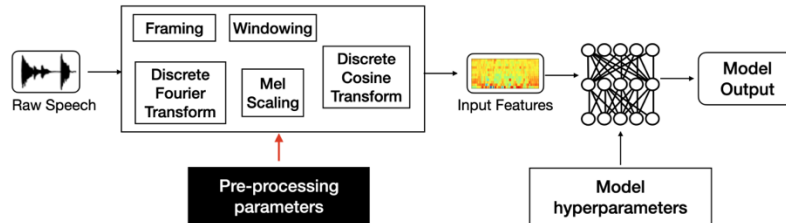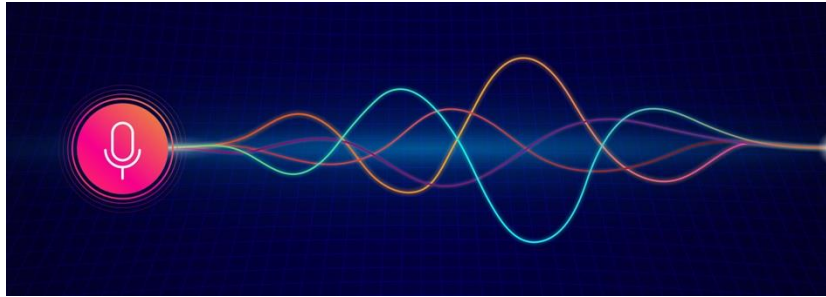  - Identify & Context
  - Algorithm & Model



Aaron Ding, Marijn Janssen, Jon Crowcroft
"Trustworthy and Sustainable Edge AI: A Research Agenda"

# *Trustworthy Edge AI*

## How ?

# Academic Case Studies

- Voice-activated services
- Vehicular services

# Case for Voice

- Poor situation…

- *"Bias exists at every development stage in the well-known VoxCeleb Speaker Recognition Challenge: model building, implementation, and data generation"*

- *"Most affected are female speakers & non-US nationalities, who experience significant performance degradation."*

"Bias Propagation in On-device ML"
**ACM TOSEM 2023**

"Bias in Automated Speaker Recognition"
**ACM FAccT 2022**

"Characterising the Role of Pre-Processing Parameters in Audio-based Embedded Machine Learning"
**ACM SenSys 2021**

# Case for Cars



- Poor situation too …

- *"Biased-car dataset leads to algorithmic bias, e.g., towards pedestrians and cyclists"*

- *"Poor data diversity… Vulnerable classes (e.g., pedestrians and cyclists) generally have less representation within the dataset"*



"Bias Detection and Generalization in AI Algorithms on Edge for Autonomous Driving"
***ACM/IEEE SEC 2022***

"Approximate Edge AI for Energy Efficient Autonomous Driving Services"
***IEEE COMST 2023***

"Adaptive Approximate Computing in Edge AI and IoT Applications"
***Elsevier JSA 2024***

**TU**Delft

# *Not Enough…*

**Real case please**

# Real Cases

Industry + Academia

## EU Horizon Project

**Score: 98/100 | Rate: 8%**

# Stakeholders



Trusted execution environment

SPATIAL
- Metrics
- XAI
- Adaptive human-readable report

AI tuning

Machine learning Subsystem
- Data
- Algorithm
- Global Model

Data

new data

feeds

Software Subsystem
- C1
- C2
- C3

feeds

AI tuning

Local Model 1
Local Model 2
......
Local Model k

Other dependecies

**Stakeholders**
- End user
- Regulating entities
- Authorities

**Stakeholders**
- Data scientist
- ML developer
- Ethics expert
- Data engineer

**Stakeholders**
- Frontend developer
- DB admin
- Backend developer
- Software tester
- Electronic engineer

**Concerns**
- Model dissected logic
- Transparency
- Accountability
- Resilience
- Safety
- Fairness

**Concerns**
- Model versioning
- Model accuracy
- Data quality
- Algorithm choice
- Privacy
- Ethics
- Framework

**Concerns**
- UI Design
- Testing
- Availability
- Database choice
- Security
- Fault tolerance
- Technology

25

# Trust as a Service

# Know-How



Domain Expertise of Consortium Partners

Literature Research

Workshops

UC1: Privacy-Preserving AI on the Edge and Beyond

UC2: Cybersecurity Analysis of 4G/5G/IoT Networks

Meetings

UC3: Accountable AI in Emergency eCall System

UC4: Resilient Cybersecurity Analytic

Requirements Analysis

Stakeholders

Data Requirements (22)

Model Requirements (14)

Legislative Requirements (15)

Security Requirements (19)

Usability Requirements (11)

Accessibility Requirements (4)

relevant for AI-based systems in general
(85 in total)

## Welcome to the Elements of AI free online course!

Join over 950,000 other people learning about the basics of AI.

## Trustworthy AI

Understand the importance, considerations, and impacts of trustworthy AI.

Buy the course

Chapter 1
**Trustworthy AI in society and business**

Chapter 2
**Fairness and accountability**

Chapter 3
**Explainability**

# Roadmap for Future



**Roadmap for Edge AI: A Dagstuhl Perspective**

Aaron Yi Ding[1]*, Ella Peltonen[2], Tobias Meuser[3], Atakan Aral[4], Christian Becker[5], Schahram Dustdar[6], Thomas Hiessl[6], Dieter Kranzlmüller[7], Madhusanka Liyanage[8], Setareh Magshudi[9], Nitinder Mohan[10], Jörg Ott[10], Jan S. Rellermeyer[11,1], Stefan Schulte[12], Henning Schulzrinne[13], Gürkan Solmaz[14], Sasu Tarkoma[15], Blesson Varghese[16], Lars Wolf[17]

ACM SIGCOMM CCR, Vol. 52, No.1, 2022

COMPUTER COMMUNICATION REVIEW

# *Sustainable Edge AI*

## How ?

# Research Agenda



Aaron Ding, Marijn Janssen, Jon Crowcroft. "Trustworthy and Sustainable Edge AI: A Research Agenda", IEEE TPS

# **Sustainable** is not a slogan

- Energy optimization for Edge AI
  - Full pipeline: data acquisition, transfer, computation, storage



Enabler

Scope

Aaron Ding, Marijn Janssen, Jon Crowcroft, "Trustworthy and Sustainable Edge AI: A Research Agenda"

# Sustainable Edge AI

- EU Marie Curie ITN: grant of €4M
- 15 Marie Curie Researchers
- 20+ industrial and academic partners

Welcome to Delft
APROPOS 2023

Approximate Computing and Adaptive Machine Learning Overlay

Distributed Cloud Computing Layer

Edge Computing Layer

Beyond 5G Smart Ecosystem

**APROPOS Project**
**Sustainable AI**

**Score: 14.5/15   Rate: 3%**

# *Case of Future Cars*

# Vehicular Data

- Data increases
- Electric cars: battery life matters!

750MB per second, as Google's driverless car prototype reported

| Autonomous Car - Sensors & Data | | |
| --- | --- | --- |
| Sensors | # Number | Data Volume |
| Camera | (8-12) | 500 - 3500 Mbit/s |
| LiDAR | (2-4) | 20-100 Mbit/s |
| Radar | (4-6) | 0.1-15 Mbit/s |
| GPS | | 50 Kb/s |
| Ultrasonic | (8-16) | 500-3500 Mb/s |
| 20 TB Car/Day | | |

Source: Elektrobit

# Energy Awareness

**Sweet spot**:

- Tolerance on accuracy & latency

- Less safety critical

- Relax the accuracy but still with acceptable experience

## How good is good enough?

# Energy Awareness

- Adaptive Approximation for

**Test-time Specialization of Dynamic Neural Networks**

IEEE CVPR MAT 2024 Best Paper Award

Sam Leroux[1,*]     Dewant Katare[2]     Aaron Yi Ding[2]     Pieter Simoens[1]

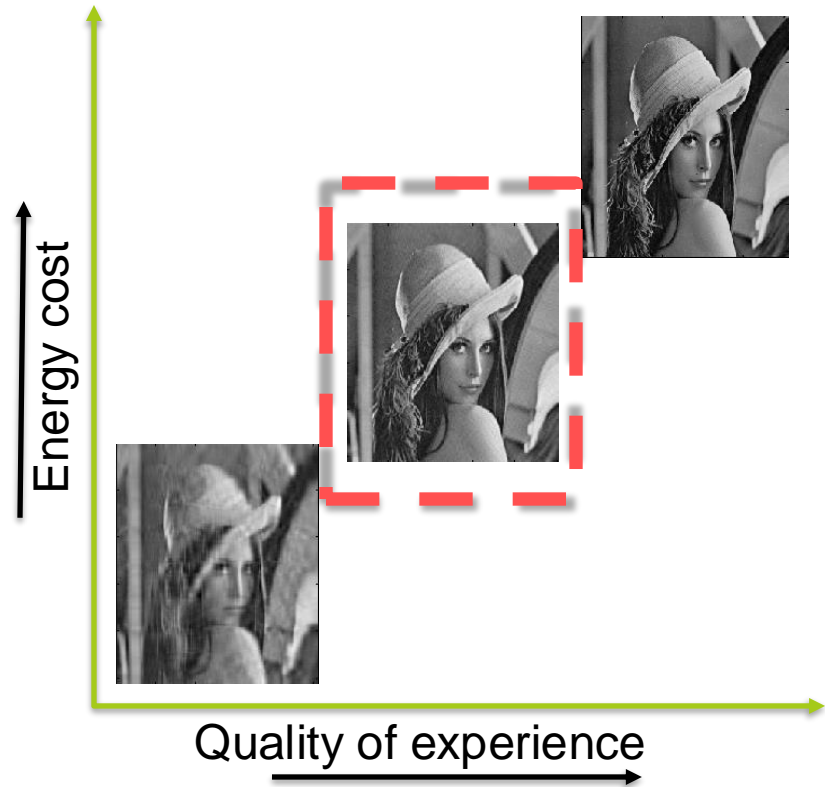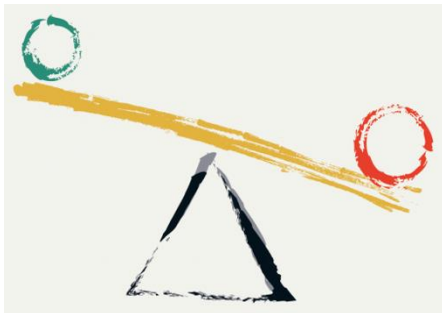[1]IDLab, Department of Information Technology, Ghent University - imec, Belgium.
[2]Department of Engineering Systems and Services, Delft University of Technology, The Netherlands.
*Corresponding author: sam.leroux@ugent.be

**Abstract**

*In recent years, there has been a notable increase in the size of commonly used image classification models. This growth has empowered models to recognize thousands of diverse object types. However, their computational demands pose significant challenges, especially when deploying them on resource-constrained edge devices. In many use cases where a model is deployed on an edge device, only a small subset of the classes will ever be observed by a given model instance. Our proposed test-time specialization of dynamic neural networks allows these models to become faster at recognizing the classes that are observed frequently, while maintaining the ability to recognize all other classes, albeit slightly less efficient. We benchmark our approach on a real-world edge device, obtaining significant speedups compared to the baseline model without test-time adaptation.*

Figure 1. We propose to first train a model on a large and diverse dataset. This model is then deployed on edge devices where it is immediately able to make useful predictions. Over time, the model is updated in a self-supervised way (test-time adaptation) to become more specialized and efficient at processing the data that is commonly observed in this environment.

"Approximate Edge AI for Energy Efficient Autonomous Driving Services"
*IEEE COMST 2023   Impact Factor 35,6*

"Nimbus: Towards Latency-Energy Efficient Task Offloading for AR Services"
*IEEE Transactions on Cloud Computing 2022*

# Independence… Grant



# Have you ever applied ?

**How many times / much** ☺

# Binary View

**MSCA**
(fellow)

**ERC**

**Erasmus+**

← Personal Grants

Project Grants →

**MSCA**
(DN)

**Horizon**
(Cluster, EIT, CHIST-ERA, Eureka ITEA-4, Eurostars-3)

Career maturity →

# Personal Grants

Academic Standing

MSCA

- **MSCA** postdoc fellowship

ERC

- **ERC** starting grant

**Erasmus+**

academic mobility

**TU**Delft

# Horizon Europe

- Time frame 2021-2027



| Pillar 1 — Excellent Science | Pillar 2 — Global Challenges and European Industrial Competitiveness | Pillar 3 — Innovative Europe |
|---|---|---|
| European Research Council | **Clusters:** Health; Culture, Creativity and Inclusive Society; Civil Security for Society; Digital, Industry and Space; Climate, Energy and Mobility; Food, Bioeconomy, Natural Resources, Agriculture and Environment | European Innovation Council |
| Marie Skłodowska-Curie Actions | | European innovation ecosystems |
| Research Infrastructures | Joint Research Centre | European Institute of Innovation and Technology |

**Widening Participation and Strengthening the European Research Area**

Widening participation and spreading excellence — Reforming and Enhancing the European R&I system

TUDelft

# Takeaway

Academic Freedom
Independence

Leadership
Industrial Connections

**MSCA**
(fellow)

**ERC**

**Erasmus+**

Personal Grants

Project Grants

**MSCA**
(DN)

**Horizon**

TU Delft

# Takeaway

## Pause Giant AI Experiments: An Open Letter

We call on all AI labs to immediately pause for at least 6 months the training of AI systems more powerful than GPT-4.

Signatures
33708

Add your signature

EDITOR: Schahram Dustdar, dustdar@dsg.tuwien.ac...

DEPARTMENT: INTERNET C...

### Revisiting Edge AI...
and Challenges

Tobias Meuser, Technical University of Dar...
Lauri Lovén, University of Oulu, 90014 Oul...
Monowar Bhuyan, Umeå University, 90187 ...
Shishir G. Patil, UC Berkeley, California, Be...
Schahram Dustdar, Vienna University of Te...
Atakan Aral, Umeå University, 90187 Umeå...
Suzan Bayhan, University of Twente, 7500 A...
Christian Becker, University of Stuttgart, 70...
Eyal de Lara, University of Toronto, Toronto...
Aaron Yi Ding, TU Delft, 2600 AA, Delft, The...
Janick Edinger, University of Hamburg, 225...
James Gross, Royal Institute of Technology...
Nitinder Mohan, Technical University of Mu...
Andy D. Pimentel, University of Amsterdam...
Etienne Rivière, UCLouvain, B-1348, Louvai...
Henning Schulzrinne, Columbia University, ...
Pieter Simoens, Ghent University-imec, B-9...
Gürkan Solmaz, NEC Laboratories Europe, 69115, Heidelberg, Germany
Michael Welzl, University of Oslo, 0313, Oslo, Norway

Edge artificial intelligence (AI) is an innovative computing paradigm that aims to shift the training and inference of machine learning models to the edge of the network. This paradigm offers the opportunity to significantly impact our everyday lives with new services such as autonomous driving and ubiquitous personalized health care. Nevertheless, bringing intelligence to the edge involves several major challenges, which include the need to constrain model architecture designs, the secure distribution and execution of the trained models, and the substantial network load required to distribute the models and data collected for training. In this article, we highlight key aspects in the development of edge AI in the past and connect them to current challenges. This article aims to identify research opportunities for edge AI, relevant to bring together the research in the fields of artificial intelligence and edge computing.

### Roadmap for Edge AI: A Dagstuhl Perspective

Aaron Yi Ding[1]*, Ella Peltonen[2], Tobias Meuser[3], Atakan Aral[4], Christian Becker[5], Schahram Dustdar[6], Thomas Hiessl[6], Dieter Kranzlmüller[7], Madhusanka Liyanage[8], Setareh Magshudi[9], Nitinder Mohan[10], Jörg Ott[10], Jan S. Rellermeyer[11,1], Stefan Schulte[12], Henning Schulzrinne[13], Gürkan Solmaz[14], Sasu Tarkoma[15], Blesson Varghese[16], Lars Wolf[17]
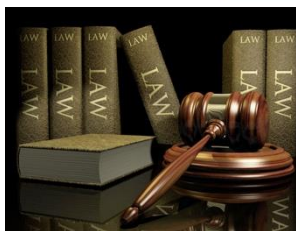
[1]TU Delft, [2]University of Oulu, [3]TU Darmstadt, [4]University of Vienna, [5]University of Mannheim, [6]TU Wien, [7]LMU Munich, [8]University College Dublin, [9]University of Tübingen, [10]TU Munich, [11]Leibniz University Hannover, [12]Hamburg University of Technology, [13]Columbia University, [14]NEC Lab, [15]University of Helsinki, [16]Queen's University Belfast, [17]TU Braunschweig

* Corresponding author: Aaron Ding (aaron.ding@tudelft.nl)

### Revisiting the Arguments for Edge Computing Research

Blesson Varghese[1], Eyal de Lara[2], Aaron Ding[3], Cheol-Ho Hong[4], Flavio Bonomi[5], Schahram Dustdar[6], Paul Harvey[7], Peter Hewkin[8], Weisong Shi[9], Mark Thiele[8], Peter Willis[10]

[1]Queen's University Belfast, UK  [2]University of Toronto, Canada  [3]TU Delft, Netherlands  [4]Chung-Ang University, S. Korea  [5]Lynx Software Technologies, USA  [6]TU Wien, Austria  [7]Rakuten Mobile, Japan  [8]SmartEdge Datacentres Ltd., UK/USA  [9]Wayne State University, USA  [10]British Telecommunications plc, UK

acm sigcomm

IEEE ComSoc
IEEE Communications Society

Personal Data
Name
Home Address
Business Address
Identity Card No
Passport No
Driving License
Income Tax No
Car Registration
Other
Confidential Data
[Identify Person]

**Enablers:** Transparency | Fairness | Robustness | Power Saving | Circular Economy | Sociotechnical Co-design

**Scoping:** Data | Identify & Context | Algorithm & Model | Trustworthy & Sustainable Edge AI | Energy | Resources | Ecosystem

**Meta Schemes:** System Design | User Engagement | Performance Tradeoff | Stakeholder | Societal Impact

"Trustworthy and Sustainable Edge AI: A Research Agenda"

TUDelft

42

# Outlook



**SEC** in Rome!

EdgeSys 2025
@ Rotterdam, NL

**ACM/IEEE Symposium on Edge Computing**



ACM SIGCOMM CCR, Vol. 52, No.1