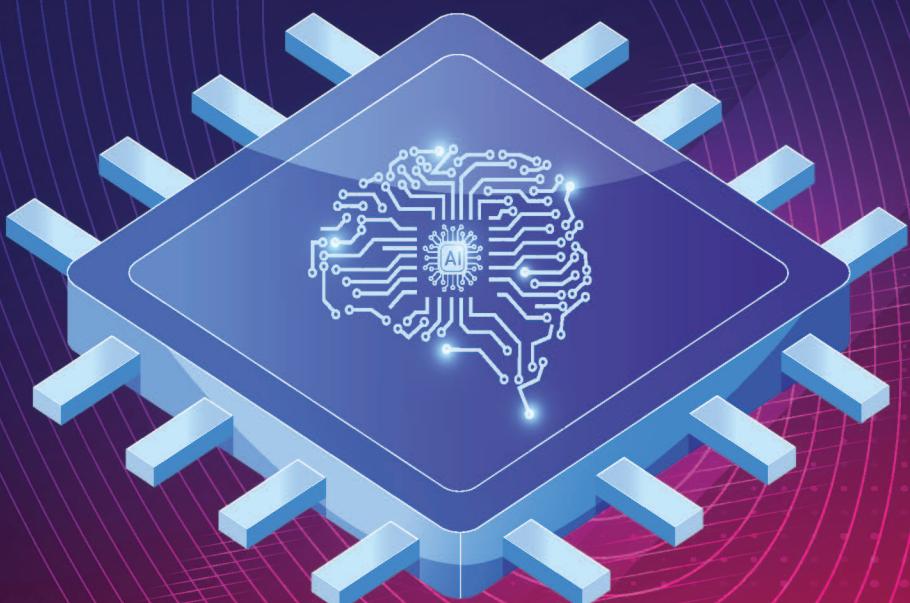


River Publishers Series in Communication and Networking

Charting the Intelligence Frontiers

Edge AI Systems Nexus

EDGE AI



Editors:

Ovidiu Vermesan

Alain Pagani

Paolo Meloni



Charting the Intelligence Frontiers Edge AI Systems Nexus

RIVER PUBLISHERS SERIES IN COMMUNICATIONS AND NETWORKING

Series Editors:

ABBAS JAMALIPOUR

*The University of Sydney
Australia*

MARINA RUGGIERI

*University of Rome Tor Vergata
Italy*

MARKO JURCEVIC

*University of Zagreb
Croatia*

The “River Publishers Series in Communications and Networking” is a series of comprehensive academic and professional books which focus on communication and network systems. Topics range from the theory and use of systems involving all terminals, computers, and information processors to wired and wireless networks and network layouts, protocols, architectures, and implementations. Also covered are developments stemming from new market demands in systems, products, and technologies such as personal communications services, multimedia systems, enterprise networks, and optical communications.

The series includes research monographs, edited volumes, handbooks and textbooks, providing professionals, researchers, educators, and advanced students in the field with an invaluable insight into the latest research and developments.

Topics included in this series include:

- Communication theory
- Multimedia systems
- Network architecture
- Optical communications
- Personal communication services
- Telecoms networks
- Wifi network protocols

For a list of other books in this series, visit www.riverpublishers.com

Charting the Intelligence Frontiers Edge AI Systems Nexus

Editors

Ovidiu Vermesan

SINTEF, Norway

Alain Pagani

German Research Center for Artificial Intelligence, Germany

Paolo Meloni

University of Cagliari, Italy



Published, sold and distributed by:

River Publishers

Broagervej 10

9260 Gistrup

Denmark

www.riverpublishers.com

ISBN: 978-87-4380-884-8 (Hardback)

978-87-4380-883-1 (Ebook)

©The Editor(s) and The Author(s) 2025. This book is published open access.

Open Access

This book is distributed under the terms of the Creative Commons Attribution-Non-Commercial 4.0 International License, CC-BY-NC 4.0 (<http://creativecommons.org/licenses/by/4.0/>), which permits use, duplication, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, a link is provided to the Creative Commons license and any changes made are indicated. The images or other third party material in this book are included in the work's Creative Commons license, unless indicated otherwise in the credit line; if such material is not included in the work's Creative Commons license and the respective action is not permitted by statutory regulation, users will need to obtain permission from the license holder to duplicate, adapt, or reproduce the material.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made.

Printed on acid-free paper.

Dedication

“Time stays long enough for those who use it.”

– Leonardo da Vinci

“Pleasure in the job puts perfection in the work.”

– Aristotle

“Success in the AI era will belong to those who adapt, learn, and innovate continuously.”

– Anonymous

“It is not the strongest species that survive, nor the most intelligent, but the ones most responsive to change.”

– Charles Darwin

Acknowledgement

The editors would like to thank all the contributors for their support in the planning and preparation of this book. The recommendations and opinions expressed in the book are those of the editors, authors, and contributors and do not necessarily represent those of any organizations, employers, or companies.

Ovidiu Vermesan
Alain Pagani
Paolo Meloni

Contents

Preface	xix
List of Figures	xxiii
List of Tables	xxix
List of Contributors	xxxi
List of Abbreviations	xxxv
1 Edge AI Systems Verification and Validation	1
<i>Ovidiu Vermesan, Alain Pagani, Roy Bahr, Marcello Antonio Coppola, and Giulio Urlini</i>	
1.1 Introduction and Background	2
1.2 Foundational Concepts and Edge AI Verification and Validation Taxonomy	6
1.2.1 Agentic AI and AI Agents	12
1.3 Defining Verification and Validation per Standard	14
1.4 Key Elements for Edge AI Verification and Validation	17
1.4.1 Core Elements for AI Verification	19
1.4.1.1 Data Verification	19
1.4.1.2 Model Verification	20
1.4.1.3 System-Level Verification	21
1.4.1.4 Process and Governance Verification	22
1.4.2 Core Elements Subject to AI Validation	23
1.4.2.1 Ensuring Fitness for Intended Purpose and Operational Context	23
1.4.2.2 Meeting User Needs and Stakeholder Expectations	23
1.4.2.3 Assessing Real-World Effectiveness and Outcomes	24

1.4.2.4	Evaluating Usability and Human-AI Interaction	24
1.4.2.5	Validating Ethical Alignment and Societal Impact	24
1.4.2.6	Data Quality and Suitability	25
1.5	The Edge AI Verification and Validation Lifecycle	26
1.6	Failure Case Behaviour in Edge-based Machine Vision Systems	29
1.7	Research Challenges in Edge AI Verification and Validation	31
1.8	Trends and Methodologies in Edge AI Verification and Validation	39
1.9	Conclusion	41
2	Pioneering the Hybridization of Federated Learning in Human Activity Recognition	53
	<i>Alfonso Esposito, Yasamin Moghbelan, Ivan Zyrianoff, Leonardo Ciabattini, Federico Montori, and Marco Di Felice</i>	
2.1	Introduction and Background	53
2.2	Hybrid FL Architecture	55
2.3	Evaluation Methodology and Metrics	57
2.4	Evaluation Results	59
2.5	Conclusion and Future Works	62
3	Edge Intelligence Architecture for Distributed and Federated Learning Systems	65
	<i>Pierluigi Dell'Acqua, Lorenzo Carnevale, and Massimo Villari</i>	
3.1	Introduction	66
3.2	Related Works	67
3.2.1	Edge Intelligence	67
3.2.2	Federated Learning	69
3.2.3	Model Compression	70
3.2.4	Beyond the State of the Art	73
3.3	Use Case	73
3.4	Architecture Proposal	75
3.4.1	Assumption	75
3.4.2	Cluster Aggregator	75
3.4.3	Cloud Components	78
3.4.4	Distributed Agent	81

3.5 Challenges	83
3.6 Conclusion	83
4 Challenges and Performance of SLAM Algorithms on Resource-constrained Devices	89
<i>Calvin Galagain, Martyna Poreba, and François Goulette</i>	
4.1 Introduction and Background	90
4.2 Related Work	91
4.3 Methodology	93
4.3.1 Selected systems	93
4.3.2 Selected systems	95
4.3.3 Evaluation metrics	96
4.3.4 Dataset	96
4.4 Experimentation	97
4.4.1 Performance evaluation	97
4.4.2 TensorRT optimization for SLAM algorithms	100
4.5 Conclusion	101
4.6 Appendix	105
4.6.1 Calculation of the metrics used in evaluation	105
4.6.1.1 Absolute Trajectory Error (ATE)	106
4.6.1.2 Relative Pose Error (RPE)	106
4.6.1.3 Frames Per Second (FPS)	107
4.6.2 Alignment methods	107
4.6.2.1 Scale alignment (scale)	108
4.6.2.2 6 Degrees of Freedom (6DOF)	108
4.6.2.3 7 Degrees of Freedom (7DOF)	109
4.6.2.4 Scale + 7 Degrees of Freedom	109
5 Designing Accelerated Edge AI Systems with Model Based Methodology	111
<i>Petri Solanti and Russell Klein</i>	
5.1 Introduction and Background	112
5.2 Model Based Cybertronic Systems Engineering	113
5.3 Designing Edge AI Systems with MBCSE Methodology	116
5.4 Creation of Bespoke AI Accelerator	118
5.4.1 High-Level Synthesis	119
5.4.2 Implementation and optimization with HLS	120
5.4.3 Verification and integration	122

5.5	Exemplary Results	122
5.6	Conclusion	123
6	Edge AI Acceleration for Critical Systems: from FPGA Hardware to CGRA Technology	127
	<i>Pietro Nannipieri, Luca Zulberti, Tommaso Pacini, Matteo Monopoli, Tommaso Bocchi, and Luca Fanucci</i>	
6.1	Introduction	127
6.2	State of the Art	129
6.3	FPG-AI: Automation Tool Flow for Efficient Deployment of Pre-trained	130
	6.3.1 Network-in-Network (NiN) case study	132
6.4	GPU@SAT: RISC-V Based SoC Featuring a Soft-GPU Hardware Accelerator	135
	6.4.1 Enhancing a soft GPU IP reliability against SEUs in space: Modelling approach and criticality analysis on a Radiation-Tolerant FPGA	136
6.5	CGR-AI: Innovative Coarse-Grained Reconfigurable Array Platform	139
6.6	Discussion and Conclusion	142
7	Model Selection and Prompting Strategies in Resource Constrained Environments for LLM-based Robotic System	147
	<i>Toms Eduards Zinars, Oskars Vismanis, Peteris Racinskis, Janis Arents, and Modris Greitans</i>	
7.1	Introduction and Background	148
7.2	Related Work	149
	7.2.1 Local Large Language Models	149
	7.2.2 Prompting	150
	7.2.3 LLM in Robotics	151
7.3	Experimental Setup	152
	7.3.1 System Description	153
	7.3.2 Testing process	154
	7.3.3 Model selection	155
	7.3.4 Testing environment	156
7.4	Results	156
	7.4.1 Differences between quantization precisions	156
	7.4.2 Differences between models	157
	7.4.3 Result comparison to VRAM usage	159

7.4.4	Result comparison to Benchmark performance	160
7.5	Conclusions	161
8	Optimising ViT for Edge Deployment: Hybrid Token Reduction for Efficient Semantic Segmentation	167
	<i>Mathilde Proust, Martyna Poreba, Calvin Galagain, Michał Szczepański, and Karim Haroun</i>	
8.1	Introduction and Background	168
8.2	Related Work	169
8.3	Methodology	170
8.3.1	Content-aware Patch Merging	172
8.3.2	Early-Pruning	173
8.4	Experiments	174
8.5	Conclusion	177
9	Recent Trends in Edge AI: Efficient Design, Training and Deployment of Machine Learning Models	181
	<i>Mark Deutel, Maen Mallah, Julio Wissing, and Stephan Scheele</i>	
9.1	Introduction	182
9.2	Scalable Deep Neural Network Architectures	182
9.2.1	Residual networks	183
9.2.2	MobileNet	184
9.2.3	EfficientNet	185
9.2.4	Scalable weights	186
9.2.5	Practical Considerations	186
9.3	Neural Architecture Search for Resource Aware DNN Deployment	187
9.3.1	Black-Box Multi-Objective optimization	188
9.3.2	Differentiable NAS	189
9.3.3	Zero-Cost neural architecture search	190
9.3.4	Practical considerations	191
9.4	Deep Neural Network Pruning	192
9.4.1	Pruning granularity	192
9.4.2	Pruning heuristics and sensitivity analysis	193
9.4.3	Magnitude or threshold based heuristics	193
9.4.3.1	L-Norm heuristics	194
9.4.3.2	Gradient Ranked Heuristics	194
9.4.3.3	Activation based heuristics	195
9.4.3.4	Relevance-based heuristics	195

9.4.4	Pruning schedule	196
9.4.5	Practical considerations	197
9.5	Quantization	197
9.5.1	Quantizers	198
9.5.2	Granularity	199
9.5.3	Methods	200
9.5.3.1	Post-Training quantization	200
9.5.3.2	Quantization-Aware training	201
9.5.4	Practical considerations	202
9.6	Cascaded Processing	203
9.6.1	Hierarchical systems	203
9.6.2	Distributed Computing	207
9.6.3	Early-Exit Neural Networks	209
9.7	Discussion	211
10	Scalable Sensor Fusion for Motion Localization in Large RF Sensing Networks	221
	<i>Fetze Pijlman</i>	
10.1	Motivation	221
10.2	Spensor Fusion via a Probabilistic Model	223
10.3	Update Equations	225
10.3.1	Update equation for $q(C_{ij} m_i = 0)$	227
10.3.2	Update equation for $q(C_{ij} m_i = 1)$	227
10.3.3	Update equation for $q(m s = 1)$	228
10.3.4	Update equation for $q(s)$	229
10.4	Conclusions and Discussion	229
11	Multi-Step Object Re-Identification on Edge Devices: A Pipeline for Vehicle Re-Identification	233
	<i>Tomass Zutis, Peteris Racinskis, Anzelika Bureka, Janis Judvaitis, Janis Arents, and Modris Greitans</i>	
11.1	Introduction	234
11.2	Related work and state of the art	235
11.2.1	Object detection	235
11.2.2	Object feature extraction	235
11.2.3	Vehicle re-identification	235
11.2.4	Available datasets	236
11.2.5	Edge implementation	237
11.3	Proposed methodology	237

11.3.1	Vehicle detection, tracking and counting	238
11.3.2	Vehicle feature extraction and storage	238
11.3.2.1	Datasets	239
11.3.2.2	Training hyper-parameters	239
11.3.3	Edge device considerations	240
11.4	Experimental settings	240
11.4.1	Receiving video from a Network camera	240
11.4.2	Vehicle re-identification	241
11.4.2.1	Testing and data annotation	241
11.4.2.2	Saving the feature extractions	241
11.5	Results	243
11.5.1	Performance metrics	243
11.5.2	Dataset generalization	243
11.5.3	Hyper-parameters	244
11.5.4	Performance on the VeRi-776 benchmark	245
11.5.5	Re-identification testing on test data from our cameras	246
11.5.5.1	Camera to camera re-identification	246
11.5.5.2	Sets of cameras	248
11.5.6	Testing the whole re-identification part of the pipeline	248
11.6	Future research	249
11.7	Conclusion	249
12	A TinyMLOps Framework for Real-world Applications	255
	<i>Mattia Antonini, Massimo Vecchio, and Fabio Antonelli</i>	
12.1	Introduction	255
12.2	TinyMLOps methodology	257
12.3	A TinyMLOps framework architecture	258
12.4	Technology Overview for TinyMLOps Adoption	261
12.5	Conclusions	263
13	Transfer and Self-learning in Probabilistic Models	267
	<i>Fetze Pijlman</i>	
13.1	Motivation	267
13.2	Prior Optimisation	269
13.3	Example Categorical Distribution	271
13.4	Conclusions and Discussion	272

14 A Novel Hierarchical Approach to Perform On-device Energy Efficient Fault Classification	275
<i>Devesh Vashishth, Julio Wissing, and Marco Wagner</i>	
14.1 Introduction and Background	276
14.2 State of the Art	277
14.2.1 Experimental setup	277
14.2.2 Related work	277
14.3 Hicnn Approach	280
14.3.1 HiCNN training	280
14.3.2 Feature forwarding	281
14.3.3 Baseline CNN and Hierarchical CNN	282
14.4 Evaluation	283
14.4.1 Experimental setup	283
14.4.2 Measurement	284
14.5 Conclusion and Future work	286
15 Discovering and Classifying Digital and Wooden Industries Products' Defects at the Edge by a Yolo/ResNet-based Approach and Beyond	289
<i>Robin Faro, Alessandro Strano, and Francesco Cancelliere</i>	
15.1 Introduction	290
15.2 Related Works	291
15.3 Spotting Defects in Wood Industry Products	294
15.3.1 Defect Detection Dataset	294
15.3.2 Experiments and Results	295
15.4 Spotting Defects in Digital Industry Products	297
15.4.1 Defect Detection and Classification Dataset	298
15.4.2 Experiments and Results	299
15.4.3 XAI Analysis: insights into ResNet-18 using Grad-CAM	300
15.5 Porting of the Models on Edge Devices	301
15.6 Conclusions and Future Works	303
16 Conscious Agents Interaction Framework for Industrial Automation	309
<i>Polina Ovsiannikova and Valeriy Vyatkin</i>	
16.1 Introduction	310
16.2 Related Research	310
16.3 Interaction Framework	314

16.3.1	Layer 1: Automation	314
16.3.2	Layer 2: Self-awareness	316
16.3.3	Layer 3: High-level communication and coordination	317
16.3.4	Layer 4: System goals	318
16.4	Case Study	319
16.4.1	Vertical farming module	319
16.4.2	HVAC system of a cruise ship	320
16.5	Discussion and Conclusion	321
17	Neuromorphic IoT Architecture for Efficient Water Management	325
<i>Mugdim Bublin, Heimo Hirner, Antoine-Martin Lanners, and Radu Grosu</i>		
17.1	Introduction and Background	326
17.2	Neuromorphic IoT Architecture	327
17.2.1	Design principles	327
17.2.2	Hierarchical distributed control and learning	327
17.3	Free Energy Principle	329
17.4	Asynchronous Processing and Event-driven Communication	330
17.5	The Role of Thresholds in Hierarchical IoT Model	331
17.5.1	Setting adaptive thresholds using prediction errors	331
17.5.2	Incorporating actions into threshold setting	332
17.5.3	Optimizing thresholds using free energy minimization	333
17.5.4	Threshold setting summary	333
17.6	Implementation	333
17.7	Case Study: Smart Village Water Management	335
17.7.1	Context and objectives	335
17.7.2	Data collection and preprocessing	336
17.7.3	Prediction models and performance	337
17.7.4	Anomaly detection	338
17.8	Discussion	339
17.8.1	Energy efficiency and communication overhead	339
17.8.2	System responsiveness and latency	339
17.8.3	Safety & security	339
17.8.4	Practical implications	339

17.9 Conclusion	340
17.10 Future Work	340
18 Online AI Benchmarking on Remote Board Farms	343
<i>Maïck Huguenin, Baptiste Dupertuis, Robin Frund, Margaux Divernois, and Nuria Pazos</i>	
18.1 Introduction and Novelty Aspect	344
18.2 State-of-the-art	345
18.3 dAIEdge-VLab Architecture	347
18.4 dAIEdge-VLab Implementation	352
18.5 Conclusion	357
19 Optimising Neural Networks for Water Stress Prediction in Europe: A Sustainable Approach	361
<i>Laura Sanz-Martín, Manal Jammal, and Javier Parra-Domínguez</i>	
19.1 Introduction and Background	362
19.2 State of the Art	364
19.3 Material and Methods	365
19.3.1 Data	365
19.3.2 Methodology	368
19.3.2.1 Data	368
19.3.2.2 Neural network architecture	368
19.3.2.3 Model optimization	369
19.3.2.4 Evaluation and metrics	372
19.4 Results	373
19.5 Conclusions	376
20 The Accountability Strikes Back: Decentralizing the Key Generation in CL-PKC with Traceable Ring Signatures	379
<i>Varesh Mishra, Aysajan Abidin, and Bart Preneel</i>	
20.1 Introduction	380
20.2 Related Work and Contributions	380
20.3 Preliminaries	381
20.3.1 Pairing Based Cryptography	381
20.3.2 Certificateless Public Key Cryptography (CL-PKC) .	382
20.3.3 Traceable Ring Signatures	382
20.3.4 Merkle Patricia Trie	383
20.4 Proposed Model	383
20.4.1 Notation	383

20.4.2	Network Architecture	384
20.4.3	Protocol	386
20.4.3.1	Network Bootstrapping	386
20.4.3.2	Key Generation	387
20.4.3.3	Full Private Key Generation	389
20.4.3.4	Update \mathcal{ST}	390
20.4.3.5	Audit Algorithm	391
20.5	Empirical Results and Analysis	396
20.6	Conclusions and Future Works	398
Index		403
About the Editors		409

Preface

A New Edge AI Reality

This book is the result of the rich exchanges of ideas and presentations at the European Conference on EDGE AI Technologies and Applications (EEAI) held on 21-23 October 2024 in Cagliari, Sardinia, Italy, offering a panoramic snapshot and a technical deep dive into the contemporary landscape of edge AI. With twenty selected chapters, it encapsulates the convergence of fundamental concepts, technical advancements, and real-world deployments that define the edge AI continuum.

Collectively, the book serves as a reference for the field, capturing the current state-of-the-art and anticipating future trends in hyperautomation, generative AI, connectivity, autonomy, and security mesh architectures. Whether you are seeking in-depth technical knowledge, inspiration for novel applications, or a strategic overview of the edge AI landscape, you will find invaluable insights from thought researchers and practitioners at the forefront of the field of edge AI.

A brief overview of each of the twenty chapters is provided below, highlighting the research and applications of edge AI that underscore the book's commitment to both technological and societal impact.

Edge AI Systems Verification and Validation: This chapter explores the challenges of verifying and validating complex edge AI systems, which integrate hardware, software, and data. It proposes a structured framework that combines model- and data-driven engineering to ensure these systems are reliable, robust, and meet regulatory standards.

Pioneering the Hybridization of Federated Learning: This work introduces a hybrid federated learning framework for human activity recognition, where some clients agree to share a portion of their data. The research assesses whether this partial data sharing can improve the overall classification accuracy of the collective model while maintaining user privacy.

Edge Intelligence Architecture for Distributed and Federated Learning: This chapter proposes a novel architecture for monitoring Electric Vehicles (EVs) by combining Federated Learning, Knowledge Distillation, and model compression. This approach enables the creation of efficient, privacy-preserving AI models that can be deployed on resource-constrained edge devices for applications like predictive maintenance.

Challenges and Performance of SLAM Algorithms on Resource-Constrained Devices: This study evaluates the performance of various visual-based SLAM (Simultaneous Localisation and Mapping) algorithms on resource-constrained hardware, such as the NVIDIA Jetson. It benchmarks several deep learning-based systems on metrics such as accuracy, energy consumption, and resource usage to assess their real-world viability.

Designing Accelerated Edge AI Systems with Model-Based Methodology: This chapter presents a Model-Based Cybertronic System Engineering (MBCSE) methodology for designing optimal edge AI systems with bespoke hardware accelerators. This approach enables a holistic analysis that balances performance, power, and cost, ensuring AI algorithms can be deployed effectively within tight system constraints.

Edge AI Acceleration for Critical Systems: Focusing on the demanding environment of satellites, this work discusses hardware solutions, such as FPGAs and CGRAs, for real-time, autonomous AI processing. The research addresses critical system challenges, including power constraints and radiation tolerance, and details the design of an FPGA-based GPU and an AI accelerator framework.

Model Selection and Prompting Strategies for LLM-Based Robotic Systems: This chapter examines the challenges of selecting and implementing Large Language Models (LLMs) in resource-constrained robotic systems. It highlights that changing model weights or precision often requires significant modifications to prompting strategies, complicating the development of modular, weight-agnostic systems.

Optimising ViT for Edge Deployment: This research presents a hybrid token reduction method, combining token merging and pruning, to make Vision Transformers (ViT) more efficient for semantic segmentation on edge devices. This approach significantly reduces computational complexity with only a minimal drop in accuracy, though it highlights challenges in exporting pruned models.

Recent Trends in Edge AI: This chapter provides a comprehensive overview of recent techniques for efficiently designing, training, and deploying machine learning models on edge devices. It covers scalable architectures, neural architecture search, and compression methods, such as quantisation and pruning, to enable energy-efficient AI in resource-limited environments.

Scalable Sensor Fusion for Motion Localization in Large RF Sensing Networks: This work addresses the challenge of accurate motion localisation in large-scale wireless sensing networks by using a probabilistic model. It demonstrates that variational Bayesian techniques offer a scalable solution for sensor fusion, enabling localised updates that model non-local effects efficiently.

Multi-Step Object Re-Identification on Edge Devices: This chapter proposes a pipeline for vehicle re-identification on edge devices using a multi-step feature extraction and matching process. The system detects an object, converts it to a vector embedding, and queries a database to find matches, achieving high precision in real-world camera network scenarios.

A TinyMLOps Framework for Real-World Applications: This work introduces a TinyMLOps framework to streamline the optimisation and deployment of AI models on microcontrollers. The framework uses cloud resources for intensive tasks while gathering real-time performance metrics from target devices, ensuring an accurate and scalable solution for deploying AI in constrained environments.

Transfer and Self-Learning in Probabilistic Models: This chapter explores the integration of transfer-learning and self-learning techniques within a single probabilistic model. The research finds that this synergy can be achieved through prior optimisation, enabling models to adapt across different environments where they are deployed.

A Novel Hierarchical Approach for On-Device Energy Efficient Fault Classification: This work proposes a hierarchical architecture utilising multiple smaller neural networks to perform energy-efficient fault classification directly on edge devices. By dividing the problem into smaller sub-tasks, the approach achieves a nine-fold reduction in energy consumption with comparable accuracy to a non-hierarchical model.

Discovering and Classifying Defects at the Edge: This chapter presents an AI-based optical inspection solution for detecting defects in digital and wooden industry products. Using YOLO and ResNet models deployed on edge

devices, the system achieves high accuracy in identifying defect positions and classifying defect types, with explainability tools clarifying the model's decisions.

Conscious Agents Interaction Framework for Industrial Automation: This paper examines the integration of human cognitive models into industrial automation, aiming to create flexible, multi-agent systems where humans and machines collaborate as equal partners. Case studies in vertical farming and HVAC control demonstrate how agents can reason and negotiate to achieve both collective and individual goals.

Neuromorphic IoT Architecture for Efficient Water Management: This work proposes a neuromorphic IoT architecture inspired by biological systems to address the energy and communication challenges of traditional IoT networks. A case study on water management demonstrates how this event-driven, asynchronous approach can be realised with neuromorphic hardware to create a more efficient and responsive system.

Online AI Benchmarking on Remote Board Farms: This project aims to create a collaborative platform, dAIEdge - VLab, that enables researchers to benchmark AI models on a range of remote edge devices. This virtual laboratory will provide access to shared resources and tools, enabling users without deep-embedded expertise to conduct live AI experiments.

Optimising Neural Networks for Water Stress Prediction in Europe: This study compares various neural network architectures and optimisers to predict water stress, a key sustainability indicator accurately. The findings show that a three-layer architecture with an Adam optimiser provides the highest accuracy, offering a valuable tool for informed water resource management.

Decentralising Key Generation in CL-PKC with Traceable Ring Signatures: This chapter addresses a key vulnerability in Federated Learning by proposing a mechanism to decentralise key generation in Certificateless Public Key Cryptography. Using traceable ring signatures and blockchain infrastructure, the model provides accountability and disincentivises malicious behaviour among trusted authorities.

List of Figures

Figure 1.1	Edge AI advantages.	3
Figure 1.2	Edge AI verification and validation process.	7
Figure 1.3	Edge AI dependability – Trustworthiness.	11
Figure 1.4	Edge AI dependability - Trustworthiness extended properties.	12
Figure 1.5	Verification and validation.	17
Figure 1.6	Verification and validation framework.	19
Figure 1.7	Edge AI system W-Model (adapted from [92]) . . .	28
Figure 2.1	Typical Federated Learning Architecture	56
Figure 2.2	Vertical Hybrid Federated Learning Architecture . .	57
Figure 2.3	Horizontal Hybrid Federated Learning Architecture	57
Figure 2.4	Implementation of the Vertical Hybridization in Flower	59
Figure 2.5	Horizontal Hybridization Results for UCI HAR . .	60
Figure 2.6	Vertical Hybridization Results for UCI HAR	61
Figure 2.7	Horizontal Hybridization Results for FEMNIST . .	61
Figure 2.8	Vertical Hybridization Results for FEMNIST	62
Figure 3.1	Six-level rating for EI described in [28].	68
Figure 3.2	In a Federated Learning scenario, each client trains its model leveraging its own private data and sends its model parameters to a central server. The central server aggregates the parameters received from each client to enhance the performance of the central global model, which is then sent back to the clients.	70
Figure 3.3	The schema illustrates the fundamental concept of KD: during the training of a simplified neural network, knowledge from a larger network is transferred to the smaller one.	72
Figure 3.4	Use case scenario.	74

Figure 3.5	Cluster Aggregator schema designed to handle FL central aggregator tasks and to implement a distillation framework adaptable during the training process.	76
Figure 3.6	Software components deployed in the Cloud.	79
Figure 3.7	The final architecture includes a Cluster Aggregator, deployed in the Cloud, and Distributed Agents, deployed on resource-constrained edge devices.	82
Figure 4.1	Overview of RDS-SLAM [1].	94
Figure 4.2	Overview of VDO-SLAM [2].	95
Figure 4.3	Trajectory predictions: each color denotes a different tested system.	98
Figure 4.4	SLAM Block Execution Time Breakdown. Left: VDO-SLAM; Right: RDS-SLAM.	100
Figure 5.1	Model-Based Cybertronics Systems Engineering Methodology	115
Figure 5.2	MBCSE Process for AI system design	116
Figure 5.3	Performance exploration	117
Figure 5.4	Alternative micro-architectures from a single source, based on tools settings and constraints	120
Figure 6.1	FPG-AI block diagram.	132
Figure 6.2	Overview of the System-on-Chip based on GPU@SAT.	136
Figure 6.3	Overview of the GPU@SAT architecture.	137
Figure 6.4	CGR-AI Engine block diagram.	140
Figure 7.1	Scematic of the MMS first demonstrator's HLP system.	153
Figure 7.2	Correctly passed tests for quantization precision comparison.	157
Figure 7.3	Correctly passed tests by various LLMs.	158
Figure 7.4	Planning success rates for the various models.	159
Figure 7.5	VRAM usage of the tested models.	160
Figure 7.6	Testing results relative to model performance on the BFCL.	160
Figure 8.1	Outline of the Proposed Hybrid Token Optimization Technique.	171
Figure 8.2	Results of Patch merging: grouped patches in blue, individual patches in red.	173

Figure 8.3	Layer-by-layer analysis considering GFLOPs and Throughput (FPS) for pruning heads placed at positions 6 and 8	176
Figure 8.4	ViT-Base segmentation results with pruned tokens masked in black	176
Figure 8.5	ViT-Tiny segmentation results with pruned tokens masked in black	177
Figure 9.1	Illustration of commonly used approaches for DNN scaling. DNNs can be scaled by either (a) widening the input of the DNN, (b) deepening the DNN by adding more layers and residual skip connections, or (c) increasing the resolution of the feature maps by adding more filters.	183
Figure 9.2	Three types of NAS are considered in the literature: (a) For black-box multi-objective optimization, many DNNs must be trained. However, optimisation results in a Pareto set of exact trade-offs between the different objectives. (b) Differentiable NAS returns only a single trade-off but is time and resource efficient because it optimizes the DNN during a single training run. (c) Zero-shot NAS allows fast DNN specialization for deployment goals, but the performance of the proposed trade-offs is only estimated.	188
Figure 9.3	The workflow of the two main quantization methods: (a) Post-Training Quantization and (b) Quantization-Aware Training (including the straight throw estimator for the backpropagation)	201
Figure 9.4	Example of simple cascade. Three classifiers are executed one after the other to classify three different labels.	203
Figure 9.5	Complex hierarchy with multiple levels. The classification flow is separated into different branches.	204
Figure 9.6	General structure of an Early-Exit CNN. After some CNN layers, exits can stop the computation based on a confidence metric	210

Figure 10.1	Downlights and indicated node pairs that are monitored for RSSI fluctuations. On the left a person working in an office and on the right a person walking on a corridor.	222
Figure 10.2	Bayesian network at a time instant showing the dependence between the various states and the sensor observations.	224
Figure 10.3	Isolated part of the network that is relevant for determining the states.	225
Figure 10.4	Example network with a mean-field assumption. . .	226
Figure 10.5	First results for two areas being named 99 and 97. The model parameters were chosen such that area 99 is a corridor while area 97 is a meeting room. In the figure one can see that isolated sensor events of 97 before 15:49 do not lead to presence in area 97 as it is more likely they originated from area 99. After probability for presence is high in area 97 then isolated sensor events of 97 at 15:50:15 do lead to an increase in presence in area 97.	230
Figure 11.1	Re-identifiable classes in smart city environments. .	234
Figure 11.2	The difference between the distribution of vehicles for the (from the left) VehicleID, VeRi and CityFlow datasets.	236
Figure 11.3	The proposed structure of the re-identification pipeline.	237
Figure 11.4	The implementation of the counting lines and Byte-Track in cameras (from the left, upper row) 1.,3. and (lower row) 2.	238
Figure 11.5	The same car visible in our network cameras (from the left) 1., 2. and 3., respectively.	240
Figure 11.6	The implementation of the saving zones (drawn as blue rectangles) as seen on the CityFlow test track video.	242
Figure 11.7	Model loss values during training on the VeRi dataset. We find values plateauing after 20 epochs. . .	244
Figure 12.1	The TinyMLOps Loop [3].	257
Figure 12.2	TinyMLOps Framework Architecture.	258

Figure 13.1	On the left-hand side a prior was sent to various diverse environments where in each environment a posterior was estimated. On the right-hand side the question is how to choose the prior for a new unknown environment when sensor events from other environments are known.	268
Figure 14.1	HiCNN distributed architecture with a divide and conquer approach of solving the classification task, with D0 node referring to no-fault class and S1,S2 and S3 nodes refers to the end of classification.	278
Figure 14.2	Baseline Architecture vs HiCNN architecture, indicating towards flexible architectures and number of layers used. The HiCNN architectures are used with different filter size and filter numbers.	282
Figure 14.3	Energy measurements setup using Raspberry pi connected to multimeter.	285
Figure 14.4	Current consumption during inference for Baseline algorithm vs Hierarchical algorithm indicating towards low latency inference by HiCNN.	286
Figure 15.1	An AOI solution consisting of an edge board for testing and a GPU server for learning.	290
Figure 15.2	Visual representation of wood defects and their distribution.	295
Figure 15.3	YOLOv8 architecture.	296
Figure 15.4	Example of prediction on validation set.	297
Figure 15.5	Training metrics of YOLOv8 model.	297
Figure 15.6	Defect-free chip and four common chip surface defects	299
Figure 15.7	Metrics' trends during training and test result. . . .	300
Figure 15.8	Metrics' trends during training and test result . . .	300
Figure 15.9	Grad-CAM activation maps for different chip surface defect classes.	301
Figure 16.1	Ontology showing the interaction of two agents. Assume-Derive phases are shown for Agent 1. . . .	314
Figure 16.2	Interaction framework.	315
Figure 16.3	Vertical farming module and its controller architecture.	319
Figure 17.1	Analogy between human nervous system and the proposed IoT architecture.	328

Figure 17.2	Black-Box view of the IoT network and the environment.	329
Figure 17.3	Memristor implementation of moving average and threshold comparison.	334
Figure 18.1	Virtual Lab Layer Model.	348
Figure 18.2	Centralized dAIEdge-VLab Architecture.	350
Figure 18.3	Decentralized dAIEdge-VLab Architecture.	351
Figure 18.4	User Input Selection through Web Interface.	354
Figure 18.5	List of Benchmarking Results.	354
Figure 18.6	Visualisation of Model Benchmarking Results . . .	355
Figure 19.1	Frecuency of water stress (%).	368
Figure 19.2	Comparison of R^2	374
Figure 19.3	Comparison of MSE and MAE.	375
Figure 20.1	The underlying architecture consists of TA nodes and Entities requesting partial private key generation. In this model, any Entity node can start the protocol with a TA node of its choice.	385
Figure 20.2	Network Bootstrapping.	386
Figure 20.3	Flowchart for the full key generation and audit procedure.	387
Figure 20.4	Key generation protocol.	388
Figure 20.5	Benchmark for traceable ring signatures for sequential and parallel execution.	396
Figure 20.6	Key generation benchmarking for $nTA = 5, 10, 20$ with different number of Entities.	397
Figure 20.7	Single run comparison of key generation for $nTA = 5, 10, 20$ with different number of Entities.	398

List of Tables

Table 3.1	From Train to Inference technologies.	72
Table 4.1	Performance Metrics: overall localization accuracy (ATE), Error between successive poses (RPE) and Inference Time	98
Table 4.2	Resource Usage	99
Table 4.3	Performance Comparison of DeepVO-pytorch Optimized with TensorRT	101
Table 6.1	Model summary of Network-in-Network.	133
Table 6.2	NiN implementation results for maximum parallelism configuration.	134
Table 6.3	Reliability values for the GPU and its components over a 60-day span.	138
Table 6.4	Classification based on criticality, area, and power. . .	139
Table 6.5	Operational & Memory Components Classes.	139
Table 6.6	Throughput compared to Area and Energy Efficiencies.	142
Table 7.1	Model Selection (all models using GGUF type Q4_K_M, all weights used in tests sourced from [46])	158
Table 8.1	Performance of Token Reduction Method integrated with ViT-Base	175
Table 8.2	Performance of Token Reduction Method integrated with ViT-Tiny	175
Table 11.1	The 3 network cameras used	241
Table 11.2	Vehicle Cropping and Saving Strategies	243
Table 11.3	Testing on our custom dataset	244
Table 11.4	Validation results during training	245
Table 11.5	VeRi Trained Model comparisons on the VeRi-776 benchmark	246
Table 11.6	Precision when re-identifying from a query camera to a gallery camera	247
Table 11.7	Recall when re-identifying from a query camera to a gallery camera	247

Table 11.8	Macro, micro precision and recall when re-identifying from a query camera to a gallery of two cameras	248
Table 11.9	Testing methods of the re-identification pipeline on CityFlow video tracks	249
Table 14.1	Comparison between Baseline and HiCNN	284
Table 15.1	Distribution of the dataset	299
Table 15.2	Performance of Nvidia Orin Nano for model deployment	302
Table 15.3	Performance of Nvidia Orin AGX for model deployment	302
Table 17.1	Input and output information available at each IoT layer.	336
Table 17.2	Hourly and daily prediction errors of different algorithms.	337
Table 19.1	Descriptive statistics for the two variables under study	366
Table 19.2	Descriptive statistics for the level of water stress of the countries under study	367
Table 19.3	Results of the different architectures and optimizers	373
Table 20.1	Notations.	385

List of Contributors

Abidin, Aysajan, COSIC, KU Leuven, Belgium

Antonelli, Fabio, Fondazione Bruno Kessler, Italy

Antonini, Mattia, Fondazione Bruno Kessler, Italy

Antonio Coppola, Marcello, STMicroelectronics, France

Arents, Janis, Institute of Electronics and Computer Science (EDI), Latvia

Bahr, Roy, SINTEF AS, Norway

Bocchi, Tommaso, University of Pisa, Italy

Bublin, Mugdim, University of Applied Science FH Campus Wien, Austria

Bureka, Anzelika, Institute of Electronics and Computer Science (EDI), Latvia

Cancelliere, Francesco, University of Catania, Italy

Carnevale, Lorenzo, University of Messina, Italy

Ciabattini, Leonardo, University of Bologna, Italy

Dell'Acqua, Pierluigi, University of Messina, Italy

Deutel, Mark, Friedrich-Alexander-Universität Erlangen-Nürnberg, Germany

Di Felice, Marco, University of Bologna, Advanced Research Center for Electronic Systems, Italy

Divernois, Margaux, Haute Ecole Arc – HES-SO, Switzerland

Dupertuis, Baptiste, Haute Ecole Arc – HES-SO, Switzerland

Eduards Zinars, Toms, Institute of Electronics and Computer Science (EDI), Latvia

Esposito, Alfonso, University of Bologna, Italy

- Fanucci, Luca**, *University of Pisa, Italy*
- Faro, Robin**, *Deepsensing SRL, Italy*
- Frund, Robin**, *Haute Ecole Arc – HES-SO, Switzerland*
- Galagain, Calvin**, *Université Paris-Saclay, CEA-List, ENSTA Paris, France*
- Goulette, François**, *ENSTA Paris, France*
- Greitans, Modris**, *Institute of Electronics and Computer Science (EDI), Latvia*
- Grosu, Radu**, *Technische Universität Wien, Austria*
- Haroun, Karim**, *Université Paris-Saclay, CEA-List, Université Côte d'Azur, France*
- Hirner, Heimo**, *University of Applied Science FH Campus Wien, Austria*
- Huguenin, Maïck**, *Haute Ecole Arc – HES-SO, Switzerland*
- Jammal, Manal**, *IoT Digital Innovation Hub, Spain*
- Judvaitis, Janis**, *Institute of Electronics and Computer Science (EDI), Latvia*
- Klein, Russell**, *Siemens EDA, USA*
- Lanners, Antoine-Martin**, *University of Applied Science FH Campus Wien, Austria*
- Mallah, Maen**, *Fraunhofer IIS, Fraunhofer Institute for Integrated Circuits, Germany*
- Mishra, Varesh**, *COSIC, KU Leuven, Belgium*
- Moghbelan, Yasamin**, *University of Bologna, Italy*
- Monopoli, Matteo**, *University of Pisa, Italy*
- Montori, Federico**, *University of Bologna, Advanced Research Center for Electronic Systems, Italy*
- Nannipieri, Pietro**, *University of Pisa, Italy*
- Ovsianikova, Polina**, *Aalto University, Finland*
- Pacini, Tommaso**, *University of Pisa, Italy*
- Pagani, Alain**, *German Research Center for Artificial Intelligence (DFKI), Germany*

- Parra-Domínguez, Javier**, *IoT Digital Innovation Hub, Spain*
- Pazos, Nuria**, *Haute Ecole Arc – HES-SO, Switzerland*
- Pijlman, Fetze**, *Signify, Eindhoven University of Technology, The Netherlands*
- Poreba, Martyna**, *Université Paris-Saclay, CEA-List, France*
- Preneel, Bart**, *COSIC, KU Leuven, Belgium*
- Proust, Mathilde**, *Université Paris-Saclay, CEA-List, France*
- Racinskis, Peteris**, *Institute of Electronics and Computer Science (EDI), Latvia*
- Sanz-Martín, Laura**, *University of Salamanca, Spain*
- Scheele, Stephan**, *Ostbayerische Technische Hochschule Regensburg, Germany*
- Solanti, Petri**, *Siemens EDA, Germany*
- Strano, Alessandro**, *Deepsensing SRL, Italy*
- Szczepanski, Michal**, *Université Paris-Saclay, CEA-List, France*
- Urlini, Giulio**, *STMicroelectronics, Italy*
- Vashishth, Devesh**, *University of Applied Sciences, Heilbronn, Germany*
- Vecchio, Massimo**, *Fondazione Bruno Kessler, Italy*
- Vermesan, Ovidiu**, *SINTEF AS, Norway*
- Villari, Massimo**, *University of Messina, Italy*
- Vismanis, Oskars**, *Institute of Electronics and Computer Science (EDI), Latvia*
- Vyatkin, Valeriy**, *Aalto University, Finland; Luleå Tekniska Universitet, Sweden*
- Wagner, Marco**, *University of Applied Sciences, Heilbronn, Germany*
- Wissing, Julio**, *Fraunhofer IIS, Fraunhofer Institute for Integrated Circuits, Germany*
- Zulberti, Luca**, *University of Pisa, Italy*
- Zutis, Tomass**, *Institute of Electronics and Computer Science (EDI), Latvia*
- Zyrianoff, Ivan**, *University of Bologna, Italy*

List of Abbreviations

AHU	Air handling unit
AI	Artificial Intelligence
ANN	Artificial Neural Network
ASIC	Application-Specific Integrated Circuit
AXI	Advanced eXtensible Interface
BDI	Belief-desire-intention
CGRA	Coarse-Grained Reconfigurable Array
CNN	Convolutional Neural Network
COTS	Commercial Off-The-Shelf
CPU	Central Processing Unit
CPU	Central Processing Unit
CTS	Content-aware Token Sharing
DL	Deep Learning
DMA	Direct Memory Access
DNN	Deep Neural Network
DSE	Design Space Exploration
DSP	Digital Signal Processing
DToP	Dynamic Token Pruning
EI	Edge Intelligence
ESA	European Space Agency
EV	Electric Vehicle
FL	Federated Learning
FP	Floating-point
FPGA	Field Programmable Gate Array
FPS	Frames Per Second
FU	Functional Unit
FPX	Fixed-point
GB	Gigabyte
GEO	Geosynchronous Earth Orbit
GPGPU	General-Purpose Computing on Graphic Processing Units

GPU	Graphics Processing Unit
HCI	Human-computer interaction
HDL	Hardware Description Language
HMI	Human-machine interface
HVAC	Heating, ventilation, and air conditioning
IoT	Internet of Things
KD	Knowledge Distillation
LEO	Low Earth Orbit
LIB	LI-ion Battery
LUT	Look Up Table
MAC	Multiply And Accumulation
MAS	Multiagent systems
MCU	Microcontroller Unit
MDE	Modular Deep Learning Engine
MES	Manufacturing execution system
mIoU	Mean Intersection Over Union
ML	Machine Learning
MM	Memory-Mapped
mmseg	MMSegmentation, an open-source semantic segmentation toolbox
MPU	Microprocessor Unit
NN	Neural Network
NPU	Neural Processing Unit
OTA	Over-the-air
PA	Power/AreaRH Radiation-Hardened
RBf	Radial Basis Functions
RHBD	Radiation-Hardened by Design
RL	Reinforcement learning
RNN	Recurrent Neural Network
RT	Radiation-Tolerant
SAN	Stochastic Activity Network
SEL	Single Event Latchup
SEU	Single Event Upset
SLAM	Simultaneous Localization and Mapping
SoC	System-on-Chip
SoC	State of Charge
SoH	State of Health
SVD	Singular Value Decomposition
TCM	Tightly-Coupled Memory

TMR	Triple Modular Redundancy
VIO	Visual-Inertial Odometry
ViT	Vision Transformer
VO	Visual Odometry
VPU	Vision Processing Unit
VRAM	video random-access memory

1

Edge AI Systems Verification and Validation

Ovidiu Vermesan¹, Alain Pagani², Roy Bahr¹,
Marcello Antonio Coppola³, and Giulio Urlini⁴

¹SINTEF AS, Norway

²German Research Center for Artificial Intelligence (DFKI), Germany

³STMicroelectronics, France

⁴STMicroelectronics, Italy

Abstract

The integration of edge artificial intelligence (AI) into different complex systems presents unique challenges, particularly concerning their reliability, robustness, safety, and transparency. Edge AI systems must function as intended and meet regulatory and technical standards. Traditional verification and validation (V&V) methodologies, which are well-suited for conventional software (SW) and hardware (HW) systems, do not fully address the unique characteristics of edge AI-based systems that include hardware, software, elements of edge AI technology stack and data.

The chapter delves into the challenges and methodologies for edge AI verification and validation to identify the unique elements required to develop verifiable edge AI systems based on a structured verification and validation framework integrated with model- and data-driven engineering principles, assurance cases, and domain-specific requirements. It highlights the terminology and concepts for edge AI as a technology that integrates HW, SW, and edge AI technology and data while presenting the challenges of the convergence of these technologies in developing verification and validation solutions.

Keywords: edge AI, edge AI system, verification, validation, machine learning, deep learning, AI agents, agentic AI, system engineering, small language models.

2 Edge AI Systems Verification and Validation

1.1 Introduction and Background

Edge AI has become a cornerstone of innovation in various industries, driving advancements in automation, decision-making, and predictive analysis. Edge AI systems applying machine learning (ML), deep learning (DL), and data processing at the edge involving deep neural networks (DNN) present significant challenges for ensuring the reliability, safety, and effectiveness of intelligent embedded devices across the edge AI computing continuum, ranging from micro- to deep- and meta-edge. Edge AI can be either deterministic or non-deterministic, based on the typical application and design choices involved. Many edge AI applications prioritise real-time, deterministic behaviour for critical tasks, such as control algorithms. Other applications can leverage the non-deterministic nature of AI to deliver more adaptable and creative solutions as the non-deterministic nature of edge AI means it can offer different interpretations based on context. In real-time applications, edge AI systems require precise timing and consistent response times. This is demanded for tasks where milliseconds of delay can be critical. Deterministic edge AI is appropriate for applications that demand predictability and consistency, while non-deterministic approaches are advantageous for applications that require adaptability, creativity, and continuous learning. The choice between using a deterministic or non-deterministic approach finally depends on the detailed requirements of the application and the expected trade-offs among predictability, adaptability, and computational cost.

The advancement of edge AI technologies and the ubiquity of automated AI-based tools have created complex operational environments. Edge AI systems are evolving towards engineering advanced adaptive systems and require new concepts for verification and validation to address the challenging multidimensional integration of HW, SW, AI models, algorithms, datasets, and the multimodality of data.

The advantages of leveraging edge AI in many industrial applications include real-time processing, enhanced privacy and data security, reduced latency, optimised bandwidth, reliability, and scalability, as illustrated in Figure 1.1.

Edge AI technology stack combines AI and IoT with edge computing, allowing data processing and edge AI algorithm execution to occur directly on devices located at the edge of the network. By bringing AI closer to the source of data generation, edge AI enables more efficient and responsive decision-making across a wide range of applications.

AI systems, particularly those based on machine learning (ML), pose unique challenges that differ from traditional software. Unlike conventional

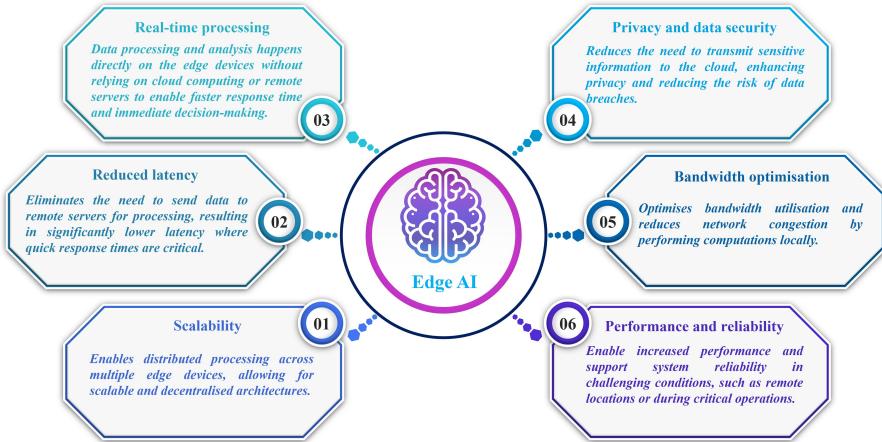


Figure 1.1 Edge AI advantages.

programs where behaviour is largely determined by explicit code, AI system behaviour often emerges from complex interactions between algorithms, vast datasets, and the operational environment [1].

Many advanced AI models, particularly deep neural networks, function as “black boxes,” making their internal decision-making processes difficult to understand, predict, or inspect directly, hence difficult to validate [2]. This opacity, combined with potential determinism/non-determinism and sensitivity to data variations, complicates efforts to guarantee reliability, safety, and fairness [3].

A particularly demanding application domain of edge AI is real-time machine vision, which is critical in domains such as industrial robotics, autonomous navigation, and quality inspection. In these systems, the correctness and timeliness of visual perception directly influence physical actions, safety, and mission success. Their dependence on high-throughput, often noisy and non-reproducible visual data, and the need for ultra-low latency, makes their verification and validation particularly challenging under edge constraints.

The data-driven approach is based on systematically and algorithmically producing the best dataset to feed a given AI-based model, focusing on improving data quality and data governance to enhance the performance of a specific problem statement. Data-driven AI aims to improve data quality and outcomes by treating code as an unchangeable entity and dealing with labelling, augmenting, managing, and curating data. This is part of the

4 Edge AI Systems Verification and Validation

data preprocessing, emphasising an iterative AI lifecycle consisting of data collection, model training, and error analysis.

The model-driven approach is based on producing the best model for a given dataset and aims to build new models and algorithmic improvements to enhance performance. The model-driven edge AI focuses on improving code reflecting the edge AI model or algorithm to achieve adequate results from fixed datasets. Edge AI developers view the training datasets from which the code, model, or algorithm is learning as a collection of reference labels. The edge AI model is made to fit that labelled training data and assumes the training data is external to the edge AI development process.

In model-driven edge AI, the focus is on optimising an edge AI model, whereas in data-driven edge AI, the focus is on data quality improvement. In model-driven edge AI, the aim is to find the most suitable edge AI model or an optimisation technique for a given problem, whereas, in data-driven edge AI, the aim is to find inconsistencies in the collected data for a given problem. The two approaches require specific verification and validation solutions.

Validation, in the context of edge AI systems, moves beyond the verification by checking if a system was built according to its technical specifications to seek confirmation that the edge AI system is fit for its intended purpose and effectively meet the actual needs and expectations of its users and stakeholders within its specific operational environment.

This necessitates the implementation of rigorous verification and validation processes, underscoring the responsibility and accountability in the development and implementation of edge AI systems.

The growing complexity and societal impact of AI, edge AI and generative AI demand a shift from purely technical verification towards a more holistic validation approach. This approach must encompass not only functional correctness but also usability, ethical alignment, fairness, robustness in real-world conditions, and overall effectiveness in achieving desired outcomes [6].

Before the adoption of AI agents and agentic AI with the use of large language models (LLMs), the development of autonomous and intelligent agents was deeply rooted in foundational paradigms of AI, such as multi-agent systems and expert systems, which emphasise social action and distributed intelligence [13][28].

Small language models (SLMs) are designed to offer capabilities similar to LLMs but scaled to edge computing capabilities, such as reduced size, processing requirements, and memory size. SLMs contain fewer parameters (e.g., hundreds of millions to one billion) while still providing strong performance for specific tasks.

Agentic AI is a class of systems that extends the capabilities of traditional AI agents by enabling multiple intelligent entities to collaborate on pursuing goals through shared memory [18][20], structured communication [24][22][26], and dynamic role assignment [21].

Ethical and legal aspects and the requirements on explainability and interpretability can lead to system development decisions that do not solely attempt to optimize functional requirements such as accuracy and robustness. In this case, system design choices rely on trade-offs that should ideally be made consciously by system developers.

Agentic AI systems pose challenges in explainability and verifiability due to their distributed, multi-agent architecture. While interpreting the behaviour of a single language model powered by the agent is already non-trivial, this complexity is multiplied when multiple agents interact asynchronously through loosely defined communication protocols. Each agent may possess its memory, task objective, and reasoning path, resulting in compounded opacity where tracing the causal chain of a final decision or failure becomes exceedingly difficult. The lack of shared, transparent logs or interpretable reasoning paths across agents makes it highly difficult, if not impossible, to determine why a particular sequence of actions occurred or which agent initiated a misstep. Compounding this opacity is the absence of formal verification tools tailored for agentic AI. In traditional software systems, model checking and formal proofs offer bounded guarantees, while there exists no widely adopted methodology to verify that a multi-agent system comprising multiple large language model agents collaborating on tasks will perform reliably across all input distributions or operational contexts.

Validation, therefore, serves as a cornerstone for building trustworthy edge AI, systems that stakeholders can confidently rely upon to operate safely, effectively, and responsibly [7]. It directly addresses the widening gap observed between accelerating edge AI capabilities and lagging safety protocols.

The AI verification standardisation efforts within the edge AI community underscores a fundamental challenge: establishing justified confidence, or trust, in edge AI systems whose behaviour often emerges unpredictably. This inherent uncertainty and the potential for significant negative impact necessitate rigorous V&V processes.

V&V encompasses activities designed to ensure that an edge AI system not only meets its specified requirements but also fulfils its intended purpose safely and reliably in its operational context. While drawing upon established V&V principles from software and systems engineering, edge AI verification

6 Edge AI Systems Verification and Validation

and validation requires tailored approaches and methodologies to address its specific complexities.

The emphasis on “trustworthiness” in standards and frameworks like ISO/IEC TR 24028 directly reflects this imperative to build demonstrable confidence in AI systems [4][5].

This chapter provides a comprehensive overview of the verification and validation of edge AI systems. It examines definitions grounded in international standards, outlines the core elements subject to verification and validation, details the typical process steps involved, analyses the significant research challenges, explores contextual variations, discusses current research trends, and summarises future directions needed to advance the field.

1.2 Foundational Concepts and Edge AI Verification and Validation Taxonomy

In edge AI systems, the failure of an AI component can lead to overall system failure, highlighting the need for AI V&V. Components with AI capabilities are treated as subsystems. V&V is carried out both on the AI subsystem itself and on its interfaces with other parts of the overall system, just as with any other subsystem. That is, the high-level definitions of V&V remain unchanged for systems containing one or more AI components.

AI V&V challenges require approaches and solutions that go beyond those for conventional or traditional systems (those without AI elements). In the context of edge AI systems, AI components and subsystems need to be integrated into the systems engineering framework. This involves identifying the characteristics of AI subsystems that create challenges in their V&V, highlighting these challenges, and providing potential solutions while determining open areas of research in the V&V of edge AI subsystems.

Conventional SW/HW systems are engineered via three main phases, namely, requirements, design and V&V. These phases are applied to each subsystem and to the system under design.

Before the expansion of AI, ML, DL, and generative AI, research on V&V of neural networks addressed the adaptation of existing standards (e.g., IEEE Std 1012-Software Verification and Validation) and processed the augmentation of these standards to enable V&V and new techniques and lessons learned to solve the V&V issues for systems integrating AI components.

In all the adaptation and augmentation attempts, one of the challenges is data validation, as the data upon which AI depends should go through a form of V&V process. Data quality attributes that are important for edge AI

systems include accuracy, currency and timeliness, correctness, consistency, usability, security and privacy, accessibility, accountability, scalability, lack of bias, and coverage and representativeness of the state space. Data validation steps can include file validation, transformation validation, import validation, domain validation, aggregation rule and business validation.

AI-based systems follow a distinct lifecycle compared to traditional systems. For edge AI systems learning lifecycle, V&V activities occur throughout the lifecycle, as illustrated in Figure 1.2. The requirements allocated to the edge AI subsystem encompass both hardware and software (HW/SW), as well as the AI models and data that flow up to the system from the edge AI subsystem.

Verification refers to the set of the activities that ensure that the edge AI system implements the specific function, and the system is built right according to requirements.

Edge AI system verification is the process of checking that the edge AI system achieves its goal without any bugs. It is the process to ensure whether the developed edge AI system is right or not. It verifies whether the developed product fulfils the requirements. Verification is static testing. Verification means answering the question: are we building the edge AI system, right?

Edge AI verification and validation require approaches and solutions at data, model and system level beyond those for cloud AI and conventional systems. Edge AI lifecycle workflows require to combine the SW/HW engineering methods with the data and system level analysis.

Data quality attributes like accuracy, timeliness, correctness, consistency, usability, security, privacy, accessibility, accountability, scalability, lack of

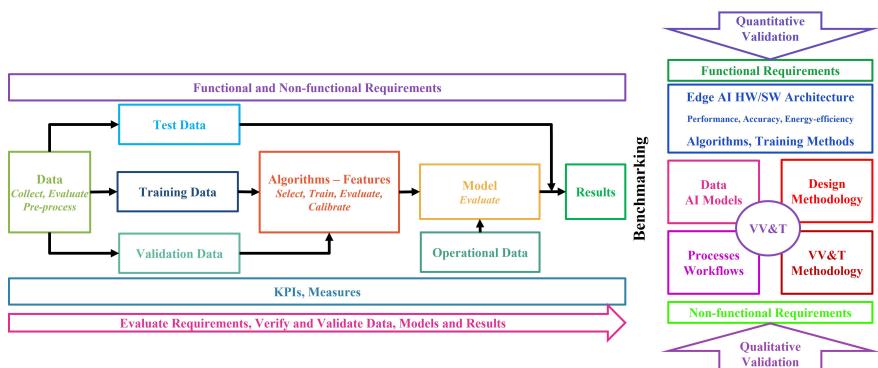


Figure 1.2 Edge AI verification and validation process.

8 Edge AI Systems Verification and Validation

bias, etc. are critical for edge AI. These data quality attributes are part of a larger edge AI non-functional requirements set.

Verification of edge AI systems involves systematically ensuring that AI models and their implementations fulfil specified requirements and intended purposes, as defined by recognised standards such as ISO/IEC 22989 (Information Technology - Artificial Intelligence - Concepts and Terminology). According to ISO, verification refers to the confirmation through objective evidence that specified requirements have been fulfilled. When applied to edge AI systems, verification processes ascertain that AI models and related software systems conform rigorously to technical and functional specifications, without necessarily validating the appropriateness of these specifications.

Principal elements involved in edge AI systems verification include:

- A formal requirements specification represents a crucial element, where clearly defined, unambiguous requirements serve as the foundational basis for verification. These specifications typically include functional requirements, performance criteria, safety constraints, security measures, and ethical guidelines.
- Model verification that entails evaluating AI models, including machine learning (ML) and deep neural networks (DNNs), ensuring their internal logic and behaviours align precisely with predefined specifications. Techniques employed in model verification include formal methods, theorem proving, model checking, and simulation-based testing.
- Software and hardware integration verification, which is vital, ensuring edge AI systems correctly interact with hardware components and software environments. It includes examining interface correctness, interoperability, real-time performance, and robustness under varying conditions and inputs.
- Rigorous test case generation and execution constitute essential verification steps. AI system verification employs automated test generation methods, including boundary value analysis, equivalence partitioning, and mutation testing, complemented by scenario-based testing to thoroughly assess compliance and performance under diverse and extreme operational conditions.
- Documentation and traceability processes involve detailed records demonstrating systematic compliance with verification steps, adherence to standards, and requirement fulfilment. Comprehensive documentation supports transparency and accountability and facilitates continuous improvement and iterative refinement processes.

In this context, the principal verification process involves several methodical steps:

- Requirement Analysis: Clearly define and document edge AI systems' functional, performance, and safety requirements.
- Verification Planning: Establishing a structured plan that details verification strategies, methods, criteria, and resources.
- Model and Code Inspection: Applying manual or automated inspections and formal verification techniques to analyse AI model structures and implementation code for correctness.
- Test Development: Generating extensive and varied test cases covering all possible usage scenarios, operational environments, and stress conditions.
- Verification Execution: Systematically conducting tests and verification activities, rigorously analysing outcomes against specified acceptance criteria.
- Reporting and Review: Documenting detailed verification outcomes, identifying discrepancies, and facilitating stakeholder review to ensure comprehensive verification coverage.
- Iterative Refinement: Addressing identified issues through iterative model adjustments, re-verification cycles, and continual improvement to achieve specified verification goals.

Validation refers to the set of the activities that ensure that the edge AI system that has been built is traceable to the requirements and the right edge AI system is built to meet user needs.

Validation is the process of checking whether the edge AI system is up to the mark or, in other words, if the product has high-level requirements. It is the process of checking the validation of the edge AI system, e.g., it checks if what we are developing is the right edge AI system. It is validation of the actual and expected edge AI systems. Validation is a form of dynamic testing. Validation means answering the question: are we building the right edge AI system?

Validation of edge AI systems is a critical and systematic process intended to ensure that the developed AI system meets stakeholders' and end-users' specific needs and expectations, as explicitly outlined in ISO/IEC 22989 (Information Technology — Artificial Intelligence — Concepts and Terminology). According to ISO standards, validation involves confirming through objective evidence that the requirements for a specific intended use or application have been fulfilled. In AI, validation goes beyond verifying

10 *Edge AI Systems Verification and Validation*

compliance with technical specifications—it assesses whether the system performs suitably in real-world conditions and scenarios.

The principal elements involved in the validation of edge AI systems encompass several dimensions:

- The identification of intended use and user requirements is foundational. Clear articulation and comprehensive understanding of user needs, operational contexts, and usage environments are paramount. This involves gathering input from stakeholders and end-users to form a robust basis for subsequent validation activities.
- Operational scenario definition is critical. Edge AI systems must be validated within scenarios that accurately represent real-world operational contexts. Scenarios are typically derived from realistic usage conditions, including normal operational states, boundary conditions, and potential abnormal or edge cases.
- Performance evaluation under realistic conditions is essential. Validation combines simulated environments and real-world testing to ensure edge AI systems perform reliably and effectively. Performance metrics, such as accuracy, precision, recall, robustness, resilience, and usability, form the basis for evaluating system performance and alignment with stakeholder expectations.
- Human-machine interaction and usability assessment are integral to validation. Edge AI systems are validated to ensure effective and intuitive interactions with human operators or users. Usability testing, user experience assessments, and feedback loops with real users facilitate comprehensive evaluations of the AI system's ease of use and accessibility.
- Safety, security, and ethical considerations are central elements of the validation process. These assessments verify that edge AI systems function correctly and comply with safety standards, security protocols, data privacy laws, and ethical guidelines, aligning with international frameworks and societal expectations.

The edge AI validation process typically involves structured, methodical steps:

- Requirement and Expectation Definition: Establishing clear validation criteria and user expectations, documenting them rigorously.
- Validation Planning: Creating detailed validation plans that specify methodologies, scenarios, test environments, and acceptance criteria.

- Scenario Development: Defining realistic operational scenarios and selecting representative use-cases and edge-cases for comprehensive validation.
- Simulation and Real-world Testing: Controlled simulations are conducted, followed by real-world trials to evaluate AI system performance against established criteria.
- Performance and Usability Assessment: Analysing performance outcomes, usability data, and user feedback to ascertain compliance with expectations and user requirements.
- Safety, Security, and Ethical Evaluation: Systematically reviewing compliance with safety and security standards, data protection requirements, and ethical norms.
- Reporting and Continuous Improvement: Compiling comprehensive validation reports, documenting findings and recommendations, and establishing iterative cycles for continuous system refinement.

Verification and validation of AI and edge AI models and data are required in safety-critical applications to ensure the trustworthiness of edge AI-enabled systems (e.g., reliability, availability, maintainability, safety, security, resilience, connectability, explainability, interpretability, transparency, etc.) as illustrated in Figure 1.3 and Figure 1.4.

Dependable edge AI systems involve using systems and software engineering principles to systematically guarantee dependability during the edge AI system's construction, V&V, and operation and consider legal and normative requirements directly from the start.

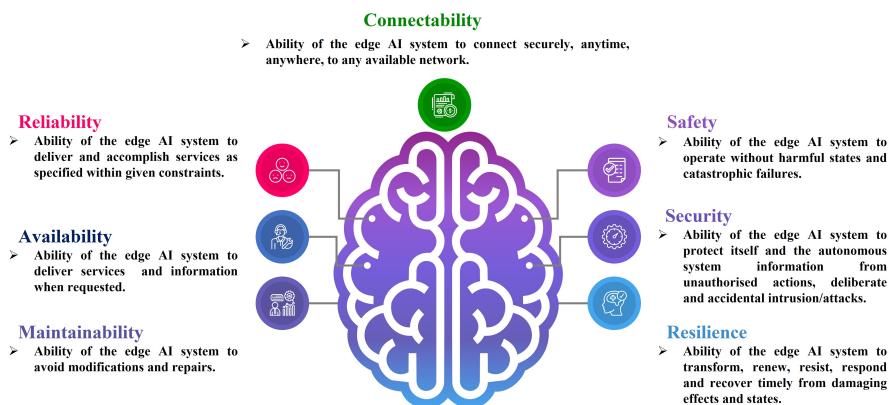


Figure 1.3 Edge AI dependability – Trustworthiness.

12 Edge AI Systems Verification and Validation

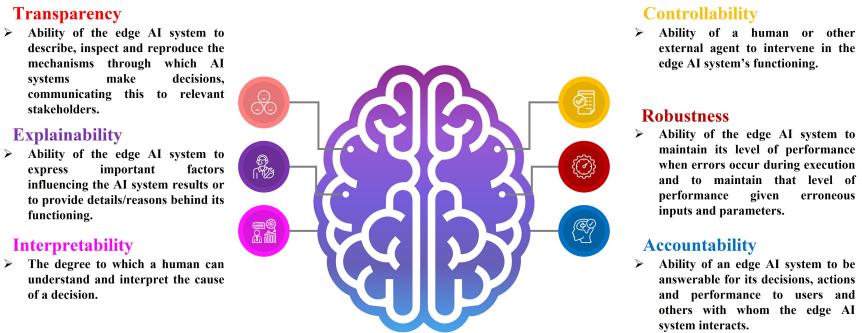


Figure 1.4 Edge AI dependability - Trustworthiness extended properties.

The progress made in developing standards and regulatory frameworks for AI and edge AI aims to ensure the responsible use of AI in various applications.

The relevant standards for AI that can be applied to edge AI systems are ISO/IEC 42001 and ISO/IEC TR 24028:2020 that are described below.

The ISO/IEC 42001 standard, a management system for AI, focuses on building trust and dependability in AI systems. It provides a framework to establish, implement, maintain, and continually improve their AI management systems, ensuring the responsible development and use of AI. The standard emphasises trustworthiness, fairness, transparency, and accountability in AI systems [43][44].

The ISO/IEC TR 24028:2020 standard addresses topics related to trustworthiness in AI systems, including approaches to establish trust in AI systems through transparency, explainability, controllability, etc.; engineering pitfalls and typical associated threats and risks to AI systems, along with possible mitigation techniques and methods; and approaches to assess and achieve availability, resiliency, reliability, accuracy, safety, security and privacy of AI systems [5].

Traditional V&V workflows, such as the V-model, are insufficient for ensuring the accuracy and reliability of AI and edge models. As a result, transformations of these workflows occurred to better serve edge AI applications.

1.2.1 Agentic AI and AI Agents

The evolution of generative AI and the emergence of AI agents and agentic AI requires addressing them under the presentation of foundational concepts

and edge AI verification and validation taxonomy by defining the concepts and their specific characteristics.

AI Agents can be defined as autonomous software entities engineered for goal-directed task execution within bounded digital environments. These agents are characterised by their ability to perceive structured or unstructured inputs, to reason over contextual information, and to initiate actions toward achieving specific objectives. The main characteristics of AI and edge AI agents are autonomy, task specificity, reactivity and adaptability, which enable the agents to operate as modular, lightweight interfaces between pre-trained AI models and domain-specific pipelines and workflows.

AI agents are the concrete instantiations of the agentic AI paradigm. An AI agent is a specific software or hardware entity that embodies the principles of agentic AI. It is a tangible system equipped with sensors to perceive its environment and effectors to act upon it. While agentic AI is the “what,” the AI agent is the “how”, the actual implementation that performs tasks, makes decisions, and interacts with external environments.

Agentic AI systems describe a paradigm shift from isolated AI agents to collaborative, multi-agent ecosystems capable of decomposing and executing complex goals [21]. These systems typically consist of orchestrated or communicating agents that interact via tools, APIs, and shared environments [23][14].

A key distinction between agentic AI and AI agents lies in their level of abstraction, as the agentic AI is a conceptual framework, whereas an AI agent is a functional system. An analogy can be drawn between the theory of computation and a physical computer. One provides the theoretical foundation and a model of what is possible, while the other is the practical machine that executes computations based on that theory.

Agentic AI reflects a broad paradigm in AI and edge AI centred on creating systems that can perceive their environment, reason about their observations, and act autonomously to achieve specific goals. It is the underlying philosophy and set of principles that guide the development of intelligent, goal-oriented systems. This concept emphasises proactivity, reactivity, and social ability, defining the potential for AI to operate as an independent actor rather than a passive tool. Agentic AI systems introduce internal orchestration mechanisms and multi-agent collaboration frameworks. Agentic AI extends the foundational architecture to support complex, distributed, and adaptive behaviours by integrating components such as specialised agents, persistent memory, orchestration and advanced reasoning and planning. Agentic AI introduces novel memory integration, communication

paradigms, and decentralised control, paving the way for the next generation of adaptive workflow automation in autonomous systems, swarm robotics, and autonomous vehicles with scalable, adaptive intelligence.

In robotics and automation, agentic AI enables collaborative behaviour in multi-robot systems. Each robot operates as a task-specialised agent, such as a picker, transporter, or mapper, while an orchestrator supervises and adapts workflows. These architectures rely on shared spatial memory, real-time sensor fusion, and inter-agent synchronisation for coordinated physical actions. Use cases include warehouse automation, drone-based orchard inspection, and robotic harvesting [25].

Verification and validation of edge AI systems, based on AI agents and agentic AI components, must focus on ensuring the correctness, reliability, and robustness of autonomous decision-making in highly dynamic and constrained environments, which requires validating that the AI agents consistently perform their intended functions correctly under varying external environment conditions, including unexpected scenarios and adversarial inputs.

Due to the limited computational resources typical of edge devices, V&V must also confirm that the system meets stringent real-time performance requirements, ensuring timely responses to critical events despite hardware and network limitations.

Another aspect is assessing the resilience and safety of adaptive learning processes within these systems, particularly as they evolve in open environments. V&V efforts should capture how individual agents and collective multi-agent behaviours emerge and interact, verifying alignment with overall system objectives and preventing unsafe or unintended actions.

Additionally, transparency and trustworthiness are key elements that enable human oversight, offering clear traceability and checkability of decisions made by autonomous components at the edge.

1.3 Defining Verification and Validation per Standard

Several ISO standards offer consistent definitions for verification and validation, primarily within the context of quality management and systems/software engineering that can be applicable to AI and edge AI as presented below.

ISO 9000:2015 (Quality management systems - Fundamentals and vocabulary) provides the definition for verification as the “confirmation,

through the provision of objective evidence, that specified requirements have been fulfilled" [29]. The focus is on confirming that the system or component conforms to its design specifications and requirements [31]. It answers the question: "Did we build the product right?" [32]. Verification is often viewed as an internal process comparing the outputs of a development phase against the inputs [31]. Validation is defined as the "confirmation, through the provision of objective evidence, that the requirements for a specific intended use or application have been fulfilled" [29]. The focus shifts to ensuring the system meets the needs of the user and fulfils its intended purpose in the actual context of use. It answers the question: "Did we build the right product?" [32]. Validation often involves testing under real or simulated use conditions and considers stakeholder needs [31].

ISO 9001:2015 (Quality management systems - Requirements), states that validation activities ensure the resulting products/services meet requirements for the specified application or intended use. Validation often involves acceptance testing with end-users and assessing fitness for purpose, making it frequently an external process, whereas verification is more often internal. Both verification and validation are essential components of quality management and are necessary for ensuring a dependable system [30].

ISO/IEC/IEEE 15288:2015 (Systems and software engineering - System life cycle processes) standard integrates V&V into the system lifecycle and considers that the **verification process** has as purpose "to provide objective evidence that a system or system element fulfils its specified requirements and characteristics" [34]. It involves activities comparing the system or element against requirements, design descriptions, and other required characteristics, confirming it was "built right" [35], while the **validation process** has as purpose "to provide objective evidence that the system, when in use, fulfils its business or mission objectives and stakeholder requirements, achieving its intended use in its intended operational environment" [33]. This process confirms that stakeholder requirements are correctly defined, and that the system meets its intended purpose in the context where it will operate [33].

ISO/IEC 22989:2022 (Information technology - Artificial intelligence - Artificial intelligence concepts and terminology) AI-specific standard defines **verification** as "confirmation, through the provision of objective evidence, that specified requirements have been fulfilled," noting it assures conformance to specification [9]. While not explicitly defining validation in the same way, it defines **trustworthiness** as the "ability to meet stakeholder

expectations in a verifiable way” [38]. This definition links the core goal of validation (meeting stakeholder expectations/needs) directly to the concept of trustworthiness in AI. The standard also incorporates a “verification and validation” phase within its depiction of the AI system lifecycle [39].

ISO/IEC TR 24028:2020 (Information technology - Artificial intelligence - Overview of trustworthiness in Artificial Intelligence) technical report further reinforces the link between validation and trustworthiness and defines **trustworthiness** as the “ability to meet stakeholder expectations in a verifiable way” [40]. This aligns the concept of trustworthiness directly with the objective of validation – confirming that stakeholder needs and intended use requirements are met [42]. The report discusses assessing and achieving key characteristics like reliability, safety, security, and privacy, all crucial aspects evaluated during validation [41].

ISO/IEC 42001:2023 (AI Management System) standard specifies requirements for establishing, implementing, maintaining, and continually improving an AI Management System (AIMS) within an organization [43]. An AIMS provides a structured framework for responsible AI governance, risk management, and operational control throughout the AI lifecycle [44]. Verification activities are integral to an AIMS, supporting risk assessment, impact assessment, performance evaluation, and ensuring compliance with policies and objectives [44]. Notably, ISO/IEC 22989 (providing the core AI terminology) is a normative reference for ISO/IEC 42001, highlighting the foundational role of clear definitions [45].

IEEE 1012-2016 (IEEE Standard for System, Software, and Hardware Verification and Validation) standard applies to systems, software, and hardware being developed, maintained, or reused (legacy, commercial off-the-shelf [COTS], non-developmental items) [91]. The term “software” also includes firmware and microcode. Additionally, each of the terms “system,” “software,” and “hardware” encompasses documentation. V&V processes include the analysis, evaluation, review, inspection, assessment, and testing of products. V&V processes are used to determine whether the development products of a given activity conform to the requirements of that activity and whether the product satisfies its intended use and user needs. V&V lifecycle process requirements are specified for different integrity levels. The scope of V&V processes encompasses systems, software, hardware, and their interfaces.

1.4 Key Elements for Edge AI Verification and Validation

The elements for verifying and validating edge AI may encompass operational aspects, system integration, AI models and human-machine interaction. Verification and validation are important to ensure reliability, performance and accuracy of complex systems.

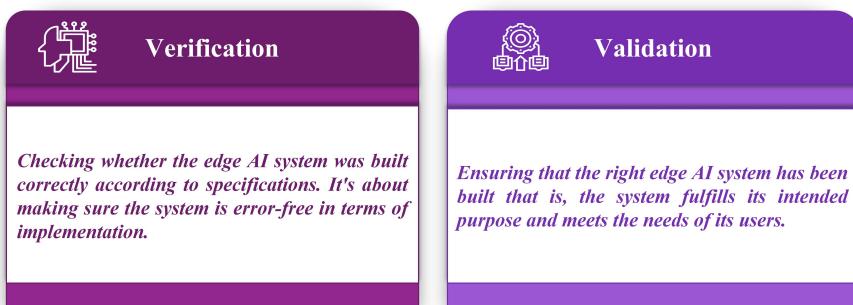


Figure 1.5 Verification and validation.

Edge AI system verification and validation refers to the processes and methodologies used to ensure that an edge AI system is dependable, performs as expected, and meets certain standards before it is deployed.

These processes are crucial as the edge AI algorithms can work with high-stakes decision-making, various sizes datasets, learn and evolve over time. The processes are needed for ensuring that the edge AI systems do what they are supposed to do, without unintended consequences, biases, or errors.

AI and edge AI systems typically focus on the actual algorithms and models to ensure that they perform as intended under various conditions. In addition, edge AI systems focus on validating the systems performance on resource-constrained devices, network conditions and privacy in real-world scenarios.

It is critical to distinguish verification from validation. While verification checks conformance to specifications (“Did we build the system right?”), validation confirms that the system meets the needs of the customer and other stakeholders and fulfils its intended purpose in its operational environment (“Did we build the right system?”) [30].

The introduction of AI in product and systems development has significantly increased the complexity of electronic components and systems (ECS), by integrating various technologies such as hardware, software, ML,

DL, NNs, generative AI, and advanced data analytics. This complexity necessitates robust verification and validation frameworks and benchmarking to ensure these systems operate correctly and efficiently as illustrated in Figure 1.6. Complex edge AI models require verification and validation to ensure their predictions, decisions, and content generation outputs are reliable and accurate, which is critical for maintaining the trustworthiness of AI systems. Failures in edge AI-based ECS can have significant economic and business-critical consequences, including system failures, financial loss, and damage to infrastructure, making the dependability of edge AI systems paramount.

In machine vision, specific verification concerns arise from the need to ensure reliable object detection, tracking, segmentation, or pose estimation across a wide range of dynamic conditions. For example, verification must confirm that visual inference results remain stable under varying lighting, occlusion, and motion blur, common challenges in edge deployments like factory floors or drones.

Ensuring robustness and reproducibility in edge-based machine vision systems is inherently difficult due to the high variability and noise in visual data. Unlike structured tabular inputs, images and videos exhibit a vast range of intra-class variation—objects or actions belonging to the same class can appear drastically different depending on factors such as:

- Lighting conditions (e.g., shadows, reflections).
- Occlusions or partial views.
- Background clutter.
- Camera distortions, blur, or motion artifacts.
- Variability in object shape, colour, texture, or viewpoint.

A comprehensive V&V framework, presented in Figure 1.6, along with benchmarking of edge AI-based methods, frameworks, tools, and ECS, is essential to ensure performance and dependable system properties like security, reliability, robustness, and fairness.

Verification ensures that edge AI-based methods, frameworks, tools, and electronic components and systems are built correctly and meet specifications, while validation confirms they perform as intended in real-world scenarios.

In edge AI systems there is a need of creating a structured approach to defining and applying such a framework to edge AI-based tools and methods, ensuring ECS meet functional and non-functional requirements, quality, KPIs, and performance standards.

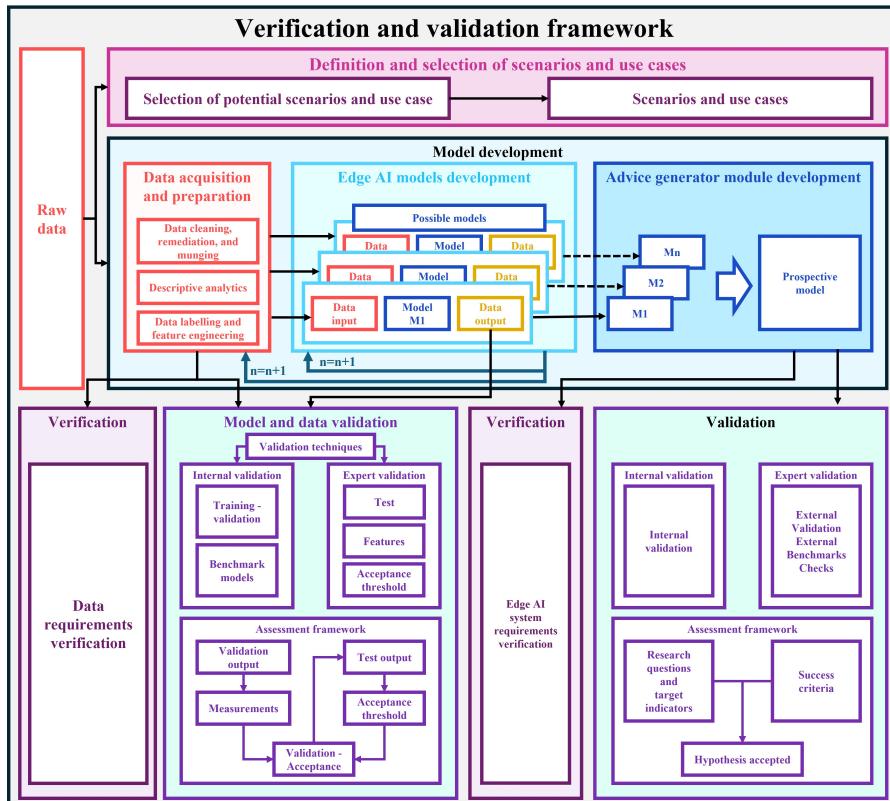


Figure 1.6 Verification and validation framework.

1.4.1 Core Elements for AI Verification

Verification activities in AI systems must address multiple facets, spanning data, models, system-level behaviour, and the processes governing development and deployment. Ensuring the integrity and appropriateness of each element is crucial for overall system trustworthiness.

1.4.1.1 Data Verification

Given that many AI and edge AI systems, particularly those based on ML, learn from data, verifying the data itself is paramount [3]. Key aspects include:

Data Quality: Assessing if the data meets predefined standards for accuracy, completeness, consistency, timeliness, and representativeness for the

target domain [3]. This involves checking for errors, missing values, correct formatting, and ensuring the data is current and relevant [46]. Poor data quality directly impacts model performance and reliability. Machine vision applications often require extensive data augmentation and synthetic dataset generation for robustness. In this case, the validity of the augmented dataset needs to be verified for plausibility and compliance with conditions of the actual use of the system.

Data Bias: Identifying systemic skews or prejudices within the data that could lead to unfair or discriminatory outcomes [3]. Verification involves confirming that bias detection methods have been applied and that any mitigation steps align with fairness requirements or definitions. This includes checking for underrepresentation or imbalances across demographic groups [3].

Data Provenance and Lineage: Ensuring the origin and history of the data are understood and documented, including all transformations and processing steps [47]. Verification confirms traceability back to authorized sources and validates the integrity of the data pipeline.

Data Security and Privacy: Confirming that data collection, storage, and processing adhere to relevant privacy regulations like General Data Protection Regulation (GDPR), a law in the European Union aimed at safeguarding the data and privacy of EU residents or California Consumer Privacy Act (CCPA) a US state law that applies to for-profit businesses operating in California that collect personal information from California residents, and organizational security policies [3]. This includes verifying the implementation of techniques like anonymization, encryption, access controls, and proper consent management [48].

Data Labelling: For supervised learning, verifying the accuracy, consistency, and quality of labels applied to the training and testing data is crucial, as errors here directly impact model learning [3].

1.4.1.2 Model Verification

The AI and edge AI model itself, the core component that performs learning and prediction, requires rigorous verification:

Accuracy and Performance: Quantifying how well the model achieves its intended task according to predefined metrics (e.g., precision, recall, F1-score for classification; BLEU score for translation) evaluated on unseen test or validation datasets [41]. Verification confirms that the achieved performance meets the specified requirements or benchmarks.

Robustness: Evaluating the model's ability to maintain its performance level when faced with noisy data, adversarial perturbations, changes in data distribution (drift), or other unexpected conditions [49]. Verification checks if the model's resilience meets specified criteria under defined stress conditions. In machine vision models, robustness testing should also include tests for perceptual artifacts, such as camera motion blur or lens distortion, and adversarial perturbations that affect visual features. This is especially important for safety-critical applications like automated visual inspection or autonomous guidance.

Reliability: Assessing the consistency and predictability of the model's outputs under normal operating conditions over time [50]. Verification aims to confirm that the model behaves dependably within its specified operational domain.

Efficiency: Measuring the model's consumption of computational resources, such as processing time, memory usage, and energy [48]. Verification ensures the model operates within the constraints imposed by the deployment hardware or system requirements.

1.4.1.3 System-Level Verification

Verification must also extend to the AI and edge AI system, considering its interaction with its environment and users:

Safety: Confirming that the system operates without causing unacceptable levels of risk or harm to humans, property, or the environment [49]. This involves verifying adherence to specific safety requirements, standards (like ISO 26262 for automotive), and risk assessments. In vision-driven systems, safety verification must ensure that the interpretation of the visual scene cannot trigger unsafe behaviour due to false positives or misclassifications e.g., mis detecting a pedestrian or failing to recognise a hazard in the camera feed.

Security and Resilience: Checking the implementation and effectiveness of measures designed to protect the system against threats like unauthorized access, data breaches, model tampering, and adversarial attacks [3]. It also includes verifying the system's ability to withstand and recover from disruptions [51].

Fairness: Evaluating system outcomes across different demographic or user groups to ensure equity and the absence of harmful bias or discrimination, according to defined fairness metrics or criteria [3].

Privacy: Verifying that the system's operation, including data handling and output generation, complies with privacy principles and regulations throughout its use [52].

1.4.1.4 Process and Governance Verification

Beyond the technical components, the processes surrounding the AI system also require verification of:

Transparency: Assessing whether sufficient and appropriate information about the AI system (its purpose, data sources, model type, limitations, performance) is documented and made available to relevant stakeholders (developers, deployers, users, regulators) [53]. Verification checks if documentation and communication channels meet specified transparency requirements.

Explainability and Interpretability: Evaluating whether the system can provide understandable reasons or justifications for its outputs or decisions, tailored to the applications and users. Verification checks if the explanation mechanisms provided meet requirements for clarity, fidelity, and utility.

Accountability: Confirming that clear roles, responsibilities, governance structures, and mechanisms for oversight, audit, and redress are defined, documented, and effectively implemented [6]. Verification involves auditing these governance processes and structures against standards like ISO/IEC 42001.

These verification elements are deeply interconnected [3]. For instance, verifying fairness requires access to appropriate data and potentially explainability techniques to understand model behaviour. Verifying safety may depend on demonstrating model robustness and having transparent documentation of system limitations. An opaque model hinders the verification of its internal logic, making it difficult to assess its safety or fairness properties directly. This interdependence necessitates a holistic verification strategy rather than treating each element in isolation.

Furthermore, the emphasis placed on different verification elements naturally shifts depending on the type of AI system. For data-driven ML models, verification heavily scrutinizes data quality, bias, model performance, and robustness [3].

In contrast, for symbolic AI systems built on explicit rules and logic, verification may concentrate more on the consistency, correctness, and completeness of the knowledge base and the soundness of the reasoning engine [64].

Hybrid neuro-symbolic systems demand verification of both the neural and symbolic parts, as well as their complex interactions, representing a distinct verification challenge [64].

1.4.2 Core Elements Subject to AI Validation

Given validation's focus on fitness for purpose and meeting stakeholder needs, the elements assessed extend beyond traditional software checks. Validating AI systems requires evaluating a broader spectrum of characteristics that reflect their performance, usability, effectiveness, and impact within their socio-technical context [6]. The exact scope may vary based on the application domain, but the core elements subject to validation include:

1.4.2.1 Ensuring Fitness for Intended Purpose and Operational Context

This is a central element of validation. It involves confirming that the AI system effectively achieves its stated goals within the specific environment and conditions of its intended use [54]. This requires a clear definition of the intended purpose and the Operational Design Domain (ODD), the specific conditions under which the system is designed to function. However, defining and validating against these can be particularly challenging for adaptive AI systems or those designed for open-world environments where conditions are dynamic and unpredictable [3]. Validation must assess performance not just under nominal conditions but also under stress, edge cases, and potential environmental shifts or adversarial inputs [85]. Frameworks like the NIST AI Risk Management Framework (RMF) emphasize establishing context (Map function) as a foundational activity to inform subsequent measurement and management, including validation [55].

1.4.2.2 Meeting User Needs and Stakeholder Expectations

Validation explicitly confirms that the system satisfies the requirements and expectations of its end-users and other relevant stakeholders [30]. This extends beyond purely functional requirements to encompass aspects like usability, user satisfaction, ease of integration into existing workflows, and alignment with business objectives [56]. Because AI systems can impact a wide range of individuals and groups, validation should involve engagement with diverse stakeholders, including end-users, domain experts, potentially affected communities, and regulators, to capture a comprehensive set of needs and expectations [90]. Addressing the challenge that these needs

might be implicit, diverse, or even conflicting is a key part of the validation process [57].

1.4.2.3 Assessing Real-World Effectiveness and Outcomes

Validation must measure how the AI and edge AI system performs in practice, assessing its actual effectiveness in achieving desired outcomes within realistic scenarios [59]. This moves beyond performance metrics derived solely from laboratory settings or curated test datasets. It involves evaluating the system's impact on relevant Key Performance Indicators (KPIs), operational efficiency, safety records, cost savings, or other context-specific measures of success [58]. Initiatives like NIST's Assessing Risks and Impacts of AI (ARIA) program are specifically focused on developing methodologies to measure these real-world impacts under controlled conditions [61]. This assessment typically requires methods such as Operational Testing (OT), field testing, pilot deployments, and continuous performance monitoring after deployment [6]. In edge machine vision systems, this includes validating that visual perception models continue to perform accurately when deployed with quantized weights, compressed inputs, or on hardware that introduces latency jitter. This real-world validation should account for degradation due to environmental variables and resource limitations.

1.4.2.4 Evaluating Usability and Human-AI Interaction

For edge AI and AI systems that interact with or support humans, validation must assess the quality and effectiveness of this interaction [60]. This includes evaluating usability (ease of use, learnability, efficiency), the clarity and utility of the interface, the cognitive load imposed on the user, and overall user satisfaction [63]. Particularly for human-AI collaboration or teaming scenarios, validation needs to assess the effectiveness of the partnership, the safety of the interaction, the appropriateness of trust levels (avoiding over-trust or under-trust), and the degree of shared understanding between human and AI [61]. This requires human-centered evaluation methods, such as usability studies, task analyses involving representative users, and systematic collection of user feedback [62].

1.4.2.5 Validating Ethical Alignment and Societal Impact

A critical dimension of edge AI validation involves assessing the system's alignment with ethical principles and societal values [6]. This includes validating characteristics like fairness, accountability, and transparency in practice [55]. Methodologies such as Ethical Impact Assessments (EIAs) are

emerging to help proactively identify, assess, and mitigate potential negative ethical and societal consequences before and during deployment [61]. A key focus is validating fairness and non-discrimination, moving beyond simple dataset metrics to assess the actual impact on different demographic groups in real-world deployment contexts [90]. This also involves considering broader societal implications related to employment, environmental sustainability, and the functioning of democratic processes [7].

1.4.2.6 Data Quality and Suitability

High-quality data ensures that models are trained effectively and can make accurate predictions in real-world scenarios [36]. As AI and edge AI systems become more complex and are deployed in diverse environments, the challenges associated with data quality and suitability have become increasingly significant. Considering the specific requirements for various AI and edge AI systems, challenges for data quality and suitability in AI and edge AI validation include:

Relevance and Representativeness: the data used for training and validation is relevant and representative of the real-world environment in which the AI and edge AI systems operate. Data must reflect the diversity of conditions, contexts, and populations that the system will encounter. If the training data is biased or unrepresentative, the model's performance may deteriorate when applied to actual situations.

Volume and Availability: Considered very important, especially in scenarios where data may be generated at high velocity. Obtaining enough high-quality data for training and validation can be difficult. In many cases, developers may struggle to gather sufficient diverse data from edge devices, leading to models that are not well-trained for all possible situations they may encounter in deployment.

Label Quality: important for supervised learning, as it directly impacts model accuracy. Inaccurate or inconsistent labelling can mislead the training process and result in poor performance in operational environments. Ensuring the reliability of labels, especially when data is labelled manually or derived from semi-automated processes, can be a significant extra work.

Bias and Fairness: the biases in learning and training of data, can lead to outcomes that are unfair when models are deployed. AI and edge AI systems trained on biased data may perpetuate existing stereotypes or discriminate against certain classes and groups. Addressing data bias and ensuring fairness

in model predictions is key to building trustworthy AI and edge systems that serve all stakeholders equitably.

Data Drift: refers to the shifts in data distributions over time, which can degrade model performance. As the underlying data evolves, models may become less accurate or irrelevant. Ongoing monitoring and adaptation of models are necessary to mitigate the effects of data drift, making it a continuous challenge for AI and edge AI validation.

Data Preprocessing: is a critical step in data management, particularly for edge AI systems with limited resources. Cleaning and transforming data into suitable formats can be challenging, when using with diverse data sources and formats. This preprocessing must be efficient to ensure real-time performance while maintaining data accuracy and integrity.

Synthetic Data: helps augment training datasets and has several limitations. The effectiveness of synthetic data depends on its ability to mimic real-world scenarios accurately. If synthetic data does not accurately represent the complexities of real-world environments, it may lead to models that underperform when applied to actual data.

Edge-Specific Challenges: these are related to data collection from distributed edge devices, considering elements like latency, bandwidth constraints, and intermittent connectivity, which can complicate the data validation process. Ensuring data quality in these scenarios requires innovative approaches to data management and model training.

1.5 The Edge AI Verification and Validation Lifecycle

According to the OECD recommendation on artificial intelligence [10], an AI system is a machine-based framework that, driven by either explicit or implicit goals, deduces from the input it receives how to produce outputs such as predictions, content, recommendations, or decisions that may impact physical or virtual environments. The levels of autonomy and adaptability of different AI systems can vary after they are deployed. The lifecycle of an AI system generally encompasses multiple stages, which include planning and design; data collection and processing; model development and/or adaptation of existing models for specific tasks; testing, evaluation, verification, and validation; deployment for use; operation and monitoring; and retirement or decommissioning. These stages often occur iteratively and are not strictly

linear. The choice to retire an AI system can be made at any time during the operation and monitoring stage.

AI and edge AI systems are distinct from other system types, which can influence the processes of the lifecycle model, such as [9]:

- Most SW systems are designed to operate in exactly defined manners dictated by their requirements and specifications. In contrast, AI and edge AI systems that utilize ML rely on data-driven training and optimization techniques to address a wide range of inputs.
- Traditional SW applications tend to be predictable, whereas this is less frequently true for AI and edge AI systems.
- Additionally, traditional SW applications are generally verifiable, while evaluating the performance of AI and edge AI systems often necessitates statistical methods, making their verification more complex.
- AI and edge AI systems usually require numerous iterations of enhancement to reach satisfactory performance levels.

The edge AI development lifecycle outlines the stages involved in creating and operationalizing edge AI systems. It starts with problem definition, functional, non-functional requirements and data collection, followed by data preparation and feature engineering. Model selection and architecture design precede the training phase, where algorithms learn from the prepared dataset. Validation and testing ensure model performance and generalization. Iterative refinement optimizes the model based on results. Deployment integrates the AI system into production environments. Monitoring and maintenance track performance, address drift, and update the model as needed.

Embedding AI and generative AI into the system design requires shifting from the current V-model, which addresses the HW and SW development cycle, to a W-model superimposed on the V-model to account for data and AI-specific artifacts. This includes the AI model development and data into the AI system's lifecycle development, as illustrated in Figure 1.7 [92].

When superimposed on the traditional V-model used in HW and SW development, the W-model AI development lifecycle creates a comprehensive framework that addresses the distinct yet interrelated processes of AI data development and HW/SW engineering. This approach ensures that AI models and supporting systems are developed in a cohesive, iterative, and validated manner. This approach aligns existing tools and methods with AI technologies. The extension into the W-model structures represents the development workflow of AI systems comprising HW, SW, AI stack, and data components.

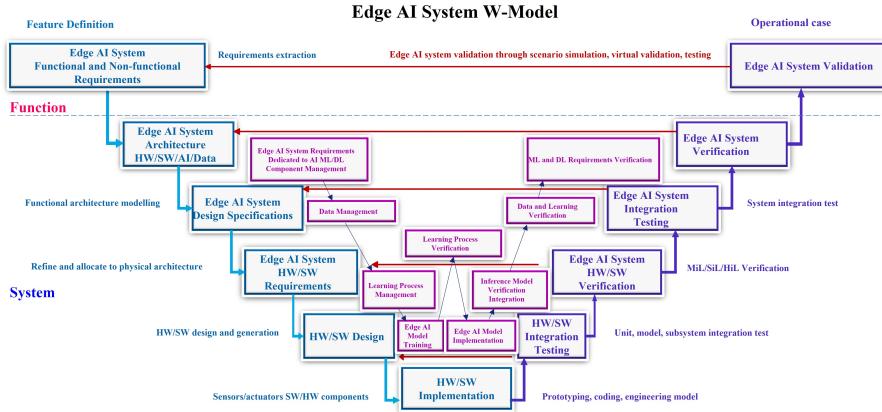


Figure 1.7 Edge AI system W-Model (adapted from [92])

The inner W part of the model represents the AI-enabled processes and workflows integrated into the conventional model. The AI system W-model emphasises systematic validation and verification at each stage of the AI development process, helping ensure the robustness, reliability, and performance of AI systems. The W-model addresses the specific design and development requirements of edge AI systems, distinguishing them from traditional HW/SW and computing paradigms. The novelty in the model lies in the fact that the data required for development and AI, ML/DL, and generative AI model training is integrated into the development cycle, superimposed on the traditional V-model, and follows the algorithm selection and training in each lifecycle stage.

The edge AI system W-model emphasises that AI and generative AI are integral to the lifecycle development processes of any AI-based product or service.

As presented in the AI system W-model at the start of the development lifecycle, developers can utilise AI and generative AI to understand domain requirements and design architecture. The design captures both functional and non-functional requirements for embedded computing systems, such as those in automotive control or industrial units, considering hardware constraints and real-time performance needs. Challenges concerning edge AI requirements and AI requirements engineering are extensive and due in part to the practice by some to treat the AI element as a “black box”. Formal specification has been attempted and has proven to be difficult for tasks that are hard to formalise, requiring decisions on the use of quantitative

or Boolean specifications, as well as the incorporation of data and formal requirements. The challenge is to design effective methods for specifying both desired and undesired properties of systems that utilise AI- or ML-based components.

When considering the broader principles of agentic AI and AI agents, their development must be integrated into the edge AI system W-model as a specific part of the lifecycle development processes of any AI-based product or service. As a result, the agentic AI and AI agents V&V extends beyond the individual agent, focusing on validating the system's autonomy and the ability to achieve long-term goals without unexpected consequences by assessing the alignment of the AI agent's goals with the overall objectives of the edge AI system and ensuring that its learning and adaptation mechanisms do not lead to unsafe or undesirable states over time.

1.6 Failure Case Behaviour in Edge-based Machine Vision Systems

In the context of edge-based machine vision systems, the study of failure case behaviour is a critical component of any robust verification and validation framework. These systems are increasingly deployed in real-world, safety-critical environments, ranging from autonomous vehicles and industrial robotics to surveillance and medical diagnostics, where failures can result in substantial consequences. While conventional validation focuses on average-case performance metrics such as accuracy or mean average precision, these metrics often obscure rare but consequential failure modes. An edge AI system that performs well under ideal conditions may fail unexpectedly in the presence of visual distortions, environmental variability, or edge hardware constraints.

Failures in machine vision models frequently arise in conditions that deviate from the data distribution seen during training. Examples include poor lighting, motion blur, occlusions, scale variation, or visual clutter. In edge deployments, such conditions are not only likely but expected, and the consequences of misclassification or missed detection can be severe. Furthermore, edge systems often operate with limited fallback options, and they must respond in real time, leaving little margin for error recovery. Understanding and characterizing these failure scenarios is therefore essential for both safety assurance and iterative model improvement.

One important strategy for investigating failure modes involves deliberate stress testing through visual perturbations. By applying controlled

transformations—such as adding noise, blurring, shifting brightness, or introducing occlusions—it becomes possible to evaluate how resilient a vision model is to real-world distortions. These tests often reveal brittle model behaviours that are not apparent during standard validation.

In addition, targeted scenario-based testing using simulation tools or recorded video sequences enables systematic exploration of edge cases. This is especially valuable for applications involving dynamic environments, such as autonomous navigation, where rare events (e.g., unexpected pedestrian appearance or sensor occlusion) may not be captured adequately in available datasets. Scenario replay or simulation also supports reproducibility of observed failures, which is often a challenge in field deployments.

Another aspect of failure case analysis is the examination of uncertainty and confidence levels in model predictions. Machine vision systems may produce incorrect predictions with unjustified confidence, especially when faced with unfamiliar or out-of-distribution inputs. Monitoring softmax confidence, prediction entropy, or Bayesian uncertainty estimates can help identify instances where the model is likely to fail. These signals may be used in runtime monitoring or to trigger fail-safe mechanisms.

Post-hoc explainability methods, such as saliency maps or activation heatmaps, also play an important role in understanding failure behaviour. By visualising the regions of an input image that contributed most to a model’s prediction, one can diagnose whether a failure was due to the model focusing on irrelevant or misleading features. This insight often reveals underlying dataset biases or spurious correlations that were inadvertently learned. By combining scenario condition variables and the predictions as features in probabilistic frameworks (e.g., Bayesian networks) is also a method for modelling the uncertainty in predictions.

Hardware-specific issues also need to be considered in failure analysis. For instance, the quantization of weights and activations required for execution on edge hardware (e.g., FPGAs or ASICs) can introduce numerical inaccuracies that degrade model performance in subtle ways. Testing the consistency between floating-point reference models and their hardware-deployed counterparts is essential to identify precision-induced errors. Similarly, real-time system profiling can reveal frame drops, synchronization mismatches, or input-output latency violations that lead to perceptual failures.

Finally, insights gained from the analysis of failure cases should feed back into the design and development process. Difficult or misclassified examples can be incorporated into retraining pipelines, synthetic data can be generated

to increase robustness, and system architectures can be adapted to detect and respond to high-uncertainty inputs. Ideally, safety envelopes are defined at design time to formally capture the operational conditions under which the system is guaranteed to function correctly. This creates a closed-loop process that not only identifies but also mitigates and prevents known failure patterns.

The systematic study of failure case behaviour is indispensable for building trustworthy machine vision systems on edge platforms. It enables developers to move beyond average-case performance toward comprehensive assurance of correctness, robustness, and safety under realistic and adverse conditions.

1.7 Research Challenges in Edge AI Verification and Validation

Edge AI represents a cutting-edge computing approach that seeks to relocate the training and inference of ML models to the network's edge [12]. However, implementing intelligence at the edge presents several significant challenges, such as the necessity to limit model architecture designs, ensuring the secure distribution and execution of the trained models, and managing the considerable network load needed to disseminate the models and the data gathered for training.

Edge AI systems that incorporate continuous learning involve the gradual updating of the models within the systems during production and test runs operations [9]. The data input into a system during these operations, is not only evaluated to generate an output but is also concurrently utilized to modify the model, aiming to enhance it based on the production data. Depending on the design of the continuous learning system, certain human interventions may be necessary, such as data labelling, validating the application of specific incremental updates, or monitoring the performance of the edge AI system. Continuous learning can address the limitations of the initial training data and assist in managing data drift and concept drift, but it also presents significant challenges in ensuring the edge AI system operates correctly while learning. It is essential to verify the system in production and to capture the production data to be able to include them as part of the training dataset in future system updates.

Catastrophic interference (and catastrophic unlearning) occurs when the training for new tasks disrupts the model's comprehension of previous tasks [11]. As new information supersedes earlier learning, the model forfeits its

capability to manage its initial tasks. Given the risk of catastrophic interference, continuous learning necessitates the capability to learn over time by integrating new observations from current data while preserving prior knowledge [9]. Numerous ML algorithms excel at learning tasks only when the data is provided in a single batch. As a model is trained on a specific task, its parameters are modified to effectively tackle that task. However, when new training data is introduced, the adjustments made for these new inputs can erase the knowledge the model had previously gained. In the context of neural networks, this occurrence is regarded as one of their key limitations.

The combination of edge AI, IoT and Cyber-Physical Systems (CPS) marks a significant transformation in data processing by bringing it closer to the origin. This strategy minimizes latency, improves real-time decision-making, and lessens the load on centralized cloud resources. In CPS, control logic is utilized to process input from sensors, through actions of actuators and thus affecting processes occurring in the physical world [9]. This is particularly evident in robotics, where sensor data is directly employed to manage the robot's operations and execute tasks in the physical world. Typically, robots are equipped with sensors at the edge to evaluate their current conditions, processors to facilitate control through analysis and action planning, and actuators to implement those actions.

In contrast to industrial robots, which are consistently repeating the same trajectories and actions without deviations, service robots or collaborative robots must adapt to evolving situations and dynamic environments [9]. Programming this adaptability presents significant challenges due to the inherent variability. Components of edge AI systems can play a role in the control software and planning processes through the “Sense-Plan-Act” framework, allowing robots to modify their actions in response to obstacles or changes in the location of target objects. The integration of robotics and edge AI system components facilitates automated physical interactions with objects, environments, and individuals.

In machine vision-based robotics, the visual processing pipeline itself must be verified and validated not only for accuracy but also for real-time responsiveness. Edge V&V must ensure that latency from image acquisition to action initiation does not exceed application-specific safety thresholds. Techniques like real-time trace logging and FPGA-based image path profiling can support this validation.

The challenges and appropriate methodologies for AI verification are not uniform; they vary significantly depending on the type of AI model employed and the application domain’s risk profile.

Model-Specific Challenges and Verification Focus

Deep Learning (DL) / Sub-Symbolic AI:

- *Challenges:* The primary verification challenges stem from their inherent opacity (making internal logic inscrutable) [64], strong data dependency (performance tied to training data quality and representativeness) [76], difficulty in formal specification of complex learned behaviours [64], susceptibility to adversarial examples, and challenges in generalization beyond training data distributions [76]. Scalability of verification methods is a major bottleneck due to the vast number of parameters and high-dimensional inputs [52]. Non-determinism can also arise during training or inference [70].
- *Verification Focus:* Emphasis is placed on empirical performance evaluation using diverse test datasets, extensive robustness testing against perturbations and adversarial attacks, fairness audits to detect biases learned from data, applying explainable AI (XAI) techniques (like LIME, SHAP, saliency maps) and interpretable AI (IAI) to gain insights into model decisions [74][75] and, where feasible, formal verification of specific, localised properties such as robustness bounds around specific inputs [69].

Symbolic AI / Rule-Based Systems:

- *Challenges:* These systems often suffer from **brittleness**, meaning they struggle to handle situations not explicitly covered by their predefined rules or knowledge base [73]. Creating and maintaining large, consistent, and complete knowledge bases can be labour-intensive and requires significant domain expertise [72]. They typically lack the ability to learn directly from raw, unstructured data. A computer vision model trained to detect stop signs may misclassify a slightly occluded or weathered sign because it hasn't seen enough variation in training. Traditional software can also exhibit brittleness, i.e. they both struggle but in different forms.
- *Verification Focus:* Verification centres on the logical integrity of the system. This includes checking the consistency of the rule set and knowledge base (absence of contradictions), analysing completeness (do the rules cover the intended domain?), formally verifying logical properties like soundness and validity of reasoning steps [64] and ensuring the traceability of outputs back to specific rules, which provides inherent explainability [72].

Neuro-symbolic AI:

- *Challenges:* This hybrid approach aims to combine the strengths of DL and symbolic AI but verifying the interaction and ensuring consistency between the neural (learning) and symbolic (reasoning) components is a key challenge [64]. Developing unified V&V frameworks that can handle both paradigms simultaneously is an active area of research [64].
- *Verification Focus:* Requires a multi-pronged approach: verifying the neural components using DL-specific techniques, verifying the symbolic components using logic-based methods, and crucially, verifying the interface and the correctness of the combined system's behaviour. A major research direction involves leveraging the symbolic part to constrain, explain, or formally verify aspects of the neural part's behaviour [64].

Domain-Specific Challenges and Verification Focus

Safety-Critical Systems (e.g., Automotive, Aerospace, Medical, Industrial Control):

- *Requirements:* These domains demand high levels of reliability, safety, robustness, and predictability [49]. System failures can have catastrophic consequences, including loss of life, severe injury, or significant environmental damage [49].
- *Challenges:* The need for provable guarantees clashes with the opacity and non-determinism of many AI components [65]. Meeting stringent regulatory standards (e.g., ISO 26262, IEC 62304, DO-178C) requires extensive evidence and documentation, which is difficult for AI/ML [68]. Managing the complexity of interaction with the physical world and ensuring safety across a vast range of operational scenarios is extremely challenging [71]. Exhaustive testing is typically infeasible due to the combinatorial explosion of possibilities [47]. Achieving deterministic replay for debugging and analysis is crucial but difficult [78].
- *Verification Focus:* Emphasis on rigorous methodologies, including formal methods where applicable, extensive simulation-based testing covering edge cases and failure modes, hardware-in-the-loop and real-world testing, fault tolerance analysis, adherence to domain-specific safety standards, meticulous documentation, and end-to-end requirements traceability [65]. Building a robust safety case with sufficient evidence is paramount [78].

Business or mission critical:

- *Requirements:* Business or mission-critical edge AI systems refer to applications that utilise AI to enable real-time decision-making and enhance operational efficiency. These domains demand high levels of scalability, reliability, safety, robustness, and interoperability. Due to their critical nature, they require special attention to ensure performance.
- *Challenges:* Deployment challenges arise from network reliability, as edge devices may operate in environments with unstable connections, affecting data synchronisation and model updates. The diversity of hardware platforms can cause compatibility issues and necessitate tailored solutions. Software challenges include the need for model optimisation, as AI models must be adjusted for edge deployment to balance accuracy and resource utilisation. Environmental conditions also pose risks, as edge devices must withstand various factors that can influence hardware performance and reliability. Regulatory and compliance challenges require navigating global data protection regulations to ensure that AI systems adhere to legal standards.
- *Verification Focus:* Verifying the effectiveness of edge AI systems involves establishing rigorous processes to ensure AI models meet performance standards under diverse conditions. Performance evaluation includes conducting real-time benchmarks to assess the responsiveness, accuracy, and resource utilisation of AI models deployed on edge devices. Interoperability ensures that edge AI solutions operate and communicate effectively within existing ecosystems and various hardware. Compliance verification requires regular audits to ensure that edge AI systems adhere to data privacy laws and industry regulations. Robustness verification involves stress-testing models against adversarial attacks and unexpected inputs to confirm their resilience in real-world scenarios. Lifecycle management strategies are necessary for overseeing the entire lifecycle of edge AI systems, from development and deployment to decommissioning.

Consumer Applications (e.g., E-commerce, social media, entertainment):

- *Requirements:* Often prioritize performance (e.g., accuracy of recommendations, speed of response), user experience, scalability, and cost-effectiveness. While direct physical safety risks are typically lower, significant concerns exist around fairness, bias, privacy, security (e.g., data breaches), misinformation, and ethical use [66].

- *Challenges:* Managing bias and fairness effectively across large, diverse user populations [66]. Protecting user privacy in data-hungry applications. Detecting and mitigating the generation or spread of harmful content or misinformation [67]. Preventing user manipulation (e.g., prompt injection in chatbots) [67]. Understanding and mitigating potential large-scale societal impacts [79][80].
- *Verification Focus:* Often relies more heavily on empirical testing, such as testing for performance, user studies for usability and acceptance, large-scale fairness and bias audits, privacy impact assessments and compliance checks, evaluation of content safety filters, and robustness testing against common failure modes or attacks. While formal verification might be used for specific critical components (e.g., payment processing), the overall verification rigor may be less intense than in safety-critical domains, unless specific high-risk functions are involved.

The fundamental difference in verification approaches between these contexts stems from the level of acceptable risk and the potential severity of failure consequences. Safety-critical domains operate with extremely low risk tolerance, demanding the highest levels of assurance and necessitating the use of more rigorous, often formal, verification techniques, alongside adherence to strict regulatory frameworks [65]. Consumer applications, while facing significant ethical and societal risks, typically have a higher tolerance for certain types of failures (e.g., a poor recommendation vs. a medical misdiagnosis), allowing for a greater reliance on empirical testing and monitoring. Addressing these domain-specific challenges in business and mission-critical edge AI systems is key for ensuring their reliability and effectiveness.

The inherent difficulties in verifying both pure DL (opacity) and pure symbolic AI (brittleness) have spurred interest in hybrid neuro-symbolic approaches [64]. By integrating the pattern-recognition strengths of neural networks with the explicit reasoning and transparency of symbolic methods, these approaches offer a potential pathway to building edge AI systems that are more amenable to verification and trust, particularly for complex tasks [77]. The verification of hybrid systems introduces its own set of research questions regarding the interaction and consistency between the different components [64].

For validation a one-size-fits-all approach to edge AI is ineffective due to the diversity of edge AI technologies and their application domains. The specific validation focus, methods, metrics, and acceptance criteria must

be tailored to the type of AI system and the context in which it operates, particularly considering the nature of user interaction and potential real-world consequences.

Validation Nuances Across AI Types

Generative AI (e.g., LLMs, SLMs, VLMs, image generators): Validation priorities include assessing factual accuracy (mitigating “hallucinations”), ensuring content safety (detecting toxicity, bias, harmful content), preventing malicious use (e.g., generating disinformation or deepfakes), and evaluating output quality attributes like coherence, relevance, and creativity, which often lack objective metrics [66]. The inherent non-determinism is a key challenge, requiring validation strategies that assess output distributions or use human evaluation and red-teaming [81]. Defining and validating the “intended purpose” for highly flexible generative models is complex [87].

Agentic AI and AI agents: The AI agent act as a deterministic component with limited scope, while agentic AI reflects distributed intelligence, characterised by goal decomposition, inter-agent communication, and contextual adaptation, demonstrating key characteristics of the modern agentic AI frameworks. Agentic AI systems define an emergent class of intelligent architectures in which multiple specialised agents collaborate to achieve complex, high-level objectives utilising collaborative reasoning and multi-step planning [17]. V&V of edge AI systems employing AI agents focuses on the reliability and safety of the agent’s actions in its operational environment to ensure the agent’s decision-making logic is robust and predictable under a broad range of inputs, especially unexpected or anomalous sensor data. This involves rigorous testing of the agent’s software, hardware, edge AI algorithms and data components to confirm they meet design specifications and performance benchmarks. Another aspect of V&V for edge AI systems that must be considered is the formal verification of the agent’s reasoning processes, which involves creating mathematical models of the agent and its environment to demonstrate that specific critical properties, such as safety, robustness, and resilience, are met. For complex, learning-based agents, this can be supplemented with extensive simulation-based testing to explore the vast state space and identify potential failure modes before deployment in the domain applications.

Autonomous Systems (e.g., autonomous vehicles, industrial robots): Validation overwhelmingly focuses on safety, reliability, and robustness

within complex and dynamic physical environments [7]. Key challenges include achieving sufficient test coverage across a vast space of potential scenarios (combinatorial explosion), bridging the gap between simulation and real-world performance, validating perception systems, and ensuring safe decision-making under uncertainty [78]. Validation heavily relies on extensive simulation, structured scenario-based testing, field testing, formal methods for safety-critical properties, and potentially runtime verification/monitoring [78]. Validating human oversight mechanisms is also critical, especially in military or safety-critical contexts [71].

Decision Support Systems (e.g., medical diagnosis aids, credit scoring tools): Validation emphasizes accuracy, reliability, fairness, explainability, and the system's impact on human decision-making and outcomes [7]. Challenges include validating against potentially imperfect or subjective ground truth, ensuring edge AI recommendations are beneficial and not misleading, rigorously assessing and mitigating bias across different user groups, and providing sufficient transparency to enable user trust and accountability. Validation typically requires domain-specific performance metrics, evaluation by domain experts, user studies assessing impact on decisions, and thorough bias and fairness audits [85].

Domain-Specific Considerations

The application domain significantly shapes validation priorities and methods due to differing risk profiles, regulatory requirements, and stakeholder concerns:

Healthcare: Extremely high stakes due to direct impact on patient safety and well-being [82]. Validation must adhere to regulatory frameworks (e.g., FDA regulations for medical devices, HIPAA for privacy, EU AI Act classifying medical AI as high-risk). Key validation elements include clinical efficacy (proven through clinical evaluation/trials), safety, reliability, data privacy, mitigation of bias in diverse patient populations, usability for clinicians, and explainability to support clinical judgment and trust. Frameworks like FUTURE-AI offer specific guidance for trustworthy AI in healthcare [86].

Finance: Focus on regulatory compliance (e.g., financial conduct authorities, anti-discrimination laws), fairness and bias mitigation in areas like credit scoring and loan applications, accuracy in fraud detection, model risk management, robustness against market volatility, security against financial attacks, and explainability for audits and customer inquiries [83]. Validation involves rigorous back testing, stress testing under various market conditions,

comprehensive bias audits using relevant fairness metrics, security penetration testing, and checks for regulatory adherence.

Transportation (especially Autonomous Vehicles): Safety is the absolute priority [82]. Validation must demonstrate safe operation under a vast range of environmental conditions (weather, lighting, road types) and interactions (other vehicles, pedestrians, cyclists). This involves validating perception systems (sensor fusion, object detection/classification), prediction models, and planning/control algorithms [71]. Validation relies heavily on extensive simulation covering millions of virtual miles, structured scenario-based testing (including edge cases and failure modes), real-world road testing, and the development of robust safety cases supported by evidence [78]. Formal verification methods may be applied to critical safety properties [71].

Social media / Content Platforms: Key concerns involve mitigating the spread of misinformation and harmful content, addressing algorithmic bias in content ranking and recommendation, ensuring fairness in content moderation, protecting user privacy, and managing the impact on user well-being and societal discourse [84]. Validation is challenging due to the massive scale, the dynamic nature of content and user behaviour, the subjectivity involved in defining “harmful” or “fair,” and the difficulty in measuring long-term societal impacts. Validation methods often include large-scale testing, human content review and rating, analysis of user engagement and feedback data, and monitoring metrics related to bias, toxicity, and content diversity.

This context-dependency highlights that effective AI validation requires not only technical expertise, but also deep domain knowledge [86]. Generic validation checklists are insufficient; protocols must be tailored to the specific AI type, its intended application, the operational environment, the relevant risks, and the specific needs and values of the stakeholders in that domain [88].

1.8 Trends and Methodologies in Edge AI Verification and Validation

The field of edge AI verification is rapidly evolving, driven by the increasing capabilities and deployment of AI systems, alongside growing concerns about their trustworthiness and potential risks. The unique challenges of edge AI validation are driving significant research and development into new methodologies, techniques, and tools. These efforts aim to provide more rigorous,

scalable, and comprehensive ways to ensure edge AI systems are fit for their intended purpose.

Formal methods are advancing with a growing interest in applying, mathematically rigorous techniques, to the verification and validation of edge AI systems. Techniques like model checking, theorem proving, abstract interpretation, and reachability analysis are being adapted to prove specific properties of edge AI components, especially neural networks, concerning safety, robustness against perturbations (e.g., adversarial examples), and fairness. Major challenges remain in scaling these methods to handle the high dimensionality and complexity of edge AI models and in formally specifying properties for systems operating under uncertainty or with incomplete requirements. Research focuses on developing more scalable algorithms, better abstraction techniques, and methods for probabilistic verification.

Explainable and interpretable AI for V&V are techniques that are increasingly explored as tools for validation and verification. By providing insights into why a model makes a certain prediction (e.g., identifying important input features using SHAP or LIME, visualizing attention mechanisms, generating counterfactual explanations), AI explainability and interpretability can help validators assess whether the model's reasoning aligns with domain knowledge, requirements, or certain rules or principles (e.g., ethical). The techniques can aid in debugging unexpected behaviours, identifying reliance on spurious correlations, and verifying compliance with constraints (e.g., fairness). This helps address the “black box” challenge for validation purposes. The reliability and interpretation of explanations themselves require validation, and research is ongoing to understand the effectiveness and limitations of using AI explainability and interpretability for V&V tasks. In the context of edge machine vision, lightweight explainability methods can help assess whether the model's attention aligns with relevant image features. These methods assist in verifying that edge vision models respond to semantically appropriate cues and not to background artifacts or compression noise.

Neuro-Symbolic AI combines the strengths of data-driven neural networks (sub-symbolic AI) with rule-based logical reasoning (symbolic AI). The symbolic component can represent explicit domain knowledge, constraints, or reasoning rules, potentially making the hybrid system more interpretable, data-efficient, and robust. From a validation perspective, neuro-symbolic approaches offer promise by potentially enabling formal verification of the symbolic reasoning part, using symbolic knowledge to constrain or validate the neural network's outputs, and providing more transparent explanations

for system behaviour. Research is actively exploring different integration architectures and their implications for validation.

Agentic AI and AI agents brings new challenges required the advancements of research focusing on developing new V&V techniques tailored to the dynamic nature of agentic AI, including advancing methods in runtime monitoring and formal verification that can cope with learning-based components and non-determinism. Creating simulation platforms that can model complex, real-world physics and multi-agent interactions will be crucial for testing edge systems exhaustively before deployment. The use of digital twin and immersive triplet environments could enable the safe exploration of an agent's behaviour under a wide range of standard and adverse conditions, helping to identify potential failure modes early. In this context, based on the technology trends research should address the system-level and collaborative aspects of agentic AI at the edge by creating frameworks for validating not only individual agents but also the collective, emergent behaviour of multi-agent systems. Developing techniques to ensure that the goals of individual agents remain aligned with the overall system objectives, even as they adapt and learn, is paramount. Research into explainable XAI and IAI for edge devices is required, as it will enable human operators to understand, trust, and effectively manage the decisions of autonomous agents, ensuring safe and predictable operation in complex, real-world scenarios.

1.9 Conclusion

The rapid advancement and deployment of edge AI necessitate a parallel evolution in the designers' ability to ensure that edge AI systems are safe, reliable, fair, and aligned with human values. Verification, as defined by standards such as ISO/IEC 22989, is the assurance through objective evidence that specified requirements have been fulfilled, forming a cornerstone of building essential trust. It provides the rigorous checks needed to confirm that AI systems are built according to their intended design and specifications.

The unique characteristics of AI, particularly its potential opacity, non-determinism, complex data dependencies, and difficulty in formally specifying requirements for emergent behaviours, pose significant challenges to traditional V&V approaches. The black-box nature of many models hinders direct inspection, scalability limitations restrict the application of formal methods, and the dynamic nature of edge AI systems and their environments demands continuous evaluation beyond design-time checks. Addressing conceptual challenges related to fairness, value alignment, and

adversarial robustness requires ongoing fundamental research, and significant progress is being made. International standards provide common terminology (ISO/IEC 22989), frameworks for trustworthiness (ISO/IEC TR 24028), and management systems for responsible AI governance (ISO/IEC 42001). Methodologically, an increasing number of researchers are adopting formal methods for specific AI and edge AI verification tasks, developing advanced testing techniques (e.g., metamorphic and adversarial testing).

Metamorphic testing techniques are used to verify the behaviour of AI models, when predicting the exact output for a given input is challenging or impossible. The metamorphic testing techniques focus on identifying relationships between inputs and outputs, known as metamorphic relations, that act as logical rules or properties that should hold true when inputs are modified. Adversarial testing is a technique in which inputs are intentionally designed to expose weaknesses or flaws in a system, thereby identifying scenarios where the system produces harmful or biased outputs. This enables testers to identify vulnerabilities and ensure the system responds safely and effectively [37].

Generative AI excels at pattern recognition, classification, and predictive analytics, generates new patterns and multimodal content (e.g., text, sound, images) and plays a dual role in the verification and validation process, for example, as part of an edge AI system that has to be verified and validated and as a technology that supports the V&V processes by generating V&V requirements, specifications and automatically performing the V&V.

The V&V of emerging edge AI agents face challenges arising from the inherent autonomy and the dynamic environments in which these agents operate. The agents can rely on machine learning models that can produce non-deterministic outputs, making the behaviour difficult to predict and formally verify. Continuous interaction with external environments introduces an extensive and unpredictable operational space, where unforeseen events can lead to emergent behaviours that were not anticipated during the design and testing phases, posing risks to safety and reliability.

Further challenges for the V&V processes are the unique constraints of the edge environment itself. Edge AI systems must operate within the limitations of computational power, memory, and energy, which can impact the performance and consistency of their decision-making processes. Edge AI agents must make real-time decisions, where latency is a critical factor. Validating that an edge AI agent responds correctly and within strict time constraints, especially when facing intermittent connectivity or degraded

sensor input, is a significant hurdle that requires novel testing methodologies beyond traditional software V&V.

The process of V&V agentic edge AI systems requires addressing the interaction between human and AI agent components to ensure that the agent's behaviour is understandable and transparent to human users, which allows for efficient oversight and human intervention when necessary. The V&V process must confirm that the edge system can communicate its state and intentions evidently and that its autonomous actions are auditable, interpretable and explainable.

This supports explainable AI techniques, implementing runtime verification for operational assurance, and opening the use of neuro-symbolic architectures to bridge the gap between learning and reasoning. Neuro-symbolic AI is a type of AI that integrates neural and symbolic AI architectures to address the weaknesses of each, providing a robust AI capable of reasoning, learning, and cognitive modelling.

Continued research is essential to develop more scalable and robust verification techniques that can handle the complexity of new edge AI systems. Addressing foundational AI safety problems, enhancing automated and human-centric V&V approaches, and building comprehensive, trustworthy AI frameworks that integrate technical verification with ethical considerations and governance are key priorities. Achieving verifiably trustworthy AI requires a holistic perspective, acknowledging the interplay between hardware, software, AI models, data, systems, processes, and the technical, application, and environmental contexts [40].

Achieving the goal of trustworthy AI and edge AI systems that are demonstrably beneficial and responsibly integrated into society requires elevating validation beyond a mere technical, end-of-phase check. It demands a holistic, continuous, and lifecycle-integrated approach [54].

This approach must rigorously integrate technical validation (ensuring robustness, reliability, and security) with user-centric validation (confirming usability, fitness for purpose, meeting needs), ethical validation (assessing fairness, accountability, and value alignment), and real-world effectiveness monitoring [90].

Success requires multidisciplinary collaboration, bringing together AI/ML experts, software engineers, domain specialists, human factors engineers, ethicists, social scientists, legal experts, end-users, and regulators. International standard bodies like ISO and organisations like NIST provide essential frameworks, common terminology, and guidance (e.g., ISO 9000,

ISO/IEC/IEEE 15288, ISO/IEC 22989, ISO/IEC TR 24028, ISO/IEC 42001, NIST AI RMF) [89].

The rapidly evolving nature of AI means that significant ongoing research and innovation in validation techniques are imperative to address the challenges effectively.

Edge AI system validation is a dynamic and increasingly critical field. As AI capabilities continue to advance and these systems become more deeply embedded in our lives, the methods used to ensure they are fit for purpose, safe, and aligned with human values must also evolve.

The focus is shifting from static, pre-deployment checks towards continuous, adaptive, and context-aware validation processes that span the entire edge AI lifecycle. Addressing the complex technical, ethical, and societal challenges associated with edge AI validation requires sustained research, multidisciplinary collaboration, and international cooperation.

Continued innovation in validation methodologies and tools will be essential to harness the transformative potential of AI responsibly and build a future where edge AI systems are trustworthy and integrated into industrial and business processes.

Acknowledgements

This publication has received funding through the projects Chips JU EdgeAI and HE dAIEDGE. The Chips JU EdgeAI “Edge AI Technologies for Optimised Performance Embedded Processing” project is supported by the Chips Joint Undertaking and its members including top-up funding by Austria, Belgium, France, Greece, Italy, Latvia, Netherlands, and Norway under grant agreement No 101097300. The HE dAIEDGE “A network of excellence for distributed, trustworthy, efficient and scalable AI at the Edge” project is supported under grant agreement No 101120726. Funded by the European Union. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the Chips Joint Undertaking. Neither the European Union nor the granting authority can be held responsible for them.

References

- [1] A. Lavin et al., “Technology readiness levels for machine learning systems,” *Nature Communications*, vol. 13, no. 1, p. 6039, Oct. 2022, <https://doi.org/10.1038/s41467-022-33128-9>.

- [2] S. Mahmud, S. Saisubramanian, and S. Zilberstein, “Verification and Validation of AI Systems Using Explanations,” *Proceedings of the AAAI Symposium Series*, vol. 4, no. 1, pp. 76–80, Nov. 2024, <https://doi.org/10.1609/aaaiiss.v4i1.31774>.
- [3] “Verification and Validation of Systems in Which AI is a Key Element - SEBoK,” *sebokwiki.org*. https://sebokwiki.org/wiki/Verification_and_Validation_of_Systems_in_Which_AI_is_a_Key_Element.
- [4] “ISO/IEC TR 24028:2020 – Overview of trustworthiness in artificial intelligence,” *BSI*, 2020. <https://www.bsigroup.com/en-IN/training-courses/isoiec-tr-240282020--overview-of-trustworthiness-in-artificial-intelligence/>.
- [5] “ISO/IEC TR 24028:2020 - Information technology - Artificial intelligence - Overview of trustworthiness in artificial intelligence” *ISO*. <https://www.iso.org/standard/77608.html>.
- [6] NIST, “AI Risk Management Framework,” *Artificial Intelligence Risk Management Framework (AI RMF 1.0)*, vol. 1, Jan. 2023, <https://doi.org/10.6028/nist.ai.100-1>.
- [7] J. Jeon, “Standardization Trends on Safety and Trustworthiness Technology for Advanced AI”, 2024, <https://arxiv.org/abs/2410.22151>.
- [8] T. R. McIntosh, T. Susnjak, T. Liu, P. Watters, and M. N. Halgamuge, “Inadequacies of Large Language Model Benchmarks in the Era of Generative Artificial Intelligence,” *arXiv (Cornell University)*, Feb. 2024. Available at: <https://doi.org/10.48550/arxiv.2402.09880>
- [9] “ISO/IEC 22989:2022 – Information technology – Artificial intelligence – Artificial intelligence concepts and terminology,” Edition 1, 2022. <https://www.iso.org/standard/74296.html>
- [10] “Recommendation of the Council on Artificial Intelligence,” OECD Legal Instruments, 2025. <https://legalinstruments.oecd.org/en/instruments/oecd-legal-0449>
- [11] “What is catastrophic forgetting?” IBM, April 2025. <https://www.ibm.com/think/topics/catastrophic-forgetting>
- [12] T. Meuser, et.al. “Revisiting Edge AI: Opportunities and Challenges”. *IEEE Internet Computing*, vol. 28, July-August 2024. <https://www.computer.org/csdl/magazine/ic/2024/04/10621659/1Z5lGDb639C>
- [13] Z. Ren and C. J. Anumba, “Multi-agent systems in construction—state of the art and prospects,” *Automation in Construction*, vol. 13, no. 3, pp. 421–434, 2004, <https://doi.org/10.1016/j.autcon.2003.12.002>.
- [14] G. Papagni, J. de Pagter, S. Zafari, M. Filzmoser, and S. T. Koeszegi, “Artificial agents’ explainability to support trust: considerations on

timing and context,” *AI & Society*, Vol. 38, No. 2, pp. 947–960, 2023. <https://doi.org/10.1007/s00146-022-01462-7>

- [15] P. Wang and H. Ding, “The rationality of explanation or human capacity? Understanding the impact of explainable artificial intelligence on human-AI trust and decision performance,” *Information Processing & Management*, Vol. 61, No. 4, p. 103732, 2024. <https://doi.org/10.1016/j.ipm.2024.103732>.
- [16] F. Sado, C. K. Loo, W. S. Liew, M. Kerzel, and S. Wermter, “Explainable Goal-driven Agents and Robots - A Comprehensive Review”, *ACM Computing Surveys*, Volume 55, Issue 10, pp. 1–41, 2023. <https://doi.org/10.1145/3564240>.
- [17] R. Sapkota, K. I. Roumeliotis, and M. Karkee, “AI Agents vs. Agentic AI: A Conceptual Taxonomy, Applications and Challenge,” *arXiv.org*, 2025. <https://arxiv.org/abs/2505.10468>.
- [18] C. Riedl and D. De Cremer, “AI for collective intelligence,” *Collective Intelligence*, vol. 4, no. 2, Apr. 2025, <https://doi.org/10.1177/26339137251328909>.
- [19] F. Piccialli, D. Chiaro, S. Sarwar, D. Cerciello, P. Qi, and V. Mele, “AgentAI: A comprehensive survey on autonomous agents in distributed AI for industry 4.0,” *Expert Systems with Applications*, vol. 291, p. 128404, Oct. 2025, <https://doi.org/10.1016/j.eswa.2025.128404>.
- [20] W. Xu, Z. Liang, K. Mei, H. Gao, J. Tan, and Y. Zhang, “A-MEM: Agentic Memory for LLM Agents,” *arXiv.org*, 2025. <https://arxiv.org/abs/2502.12110>.
- [21] D. B. Acharya, K. Kuppan and B. Divya, “Agentic AI: Autonomous Intelligence for Complex Goals—A Comprehensive Survey,” in *IEEE Access*, vol. 13, pp. 18912–18936, 2025, <https://www.doi.org/10.1109/ACCESS.2025.3532853>.
- [22] R. Zhang et al., “Toward Agentic AI: Generative Information Retrieval Inspired Intelligent Communications and Networking,” *arXiv.org*, 2025. <https://arxiv.org/abs/2502.16866>.
- [23] M. Gridach, J. Nanavati, K. Zine, L. Mendes, and C. Mack, “Agentic AI for Scientific Discovery: A Survey of Progress, Challenges, and Future Directions,” *arXiv.org*, 2025. <https://arxiv.org/abs/2503.08979>.
- [24] E. Miehling et al., “Agentic AI Needs a Systems Theory,” *arXiv.org*, 2025. <https://arxiv.org/abs/2503.00237>.
- [25] S. Hong et al., “MetaGPT: Meta Programming for Multi-Agent Collaborative Framework,” *arXiv.org*, Aug. 07, 2023. <https://arxiv.org/abs/2308.00352>.

- [26] U. M. Borghoff, P. Bottoni, and R. Pareschi, “Human-artificial interaction in the age of agentic AI: a system-theoretical approach,” *Frontiers in Human Dynamics*, vol. 7, May 2025, <https://doi.org/10.3389/fhumd.2025.1579166>.
- [27] J. Heer, “Agency plus automation: Designing artificial intelligence into interactive systems,” *Proceedings of the National Academy of Sciences*, Vol. 116, No. 6, pp. 1844–1850, 2019. <https://doi.org/10.1073/pnas.1807184115>.
- [28] E. Oliveira, K. Fischer, and O. Stepankova, “Multi-agent systems: which research for which applications,” *Robotics and Autonomous Systems*, vol. 27, no. 1-2, pp. 91–106, 1999, [https://doi.org/10.1016/S0921-8890\(98\)00085-2](https://doi.org/10.1016/S0921-8890(98)00085-2).
- [29] Validation - Glossary | CSRC - NIST Computer Security Resource Center, <https://csrc.nist.gov/glossary/term/validation>
- [30] Verification and validation - Wikipedia, https://en.wikipedia.org/wiki/Verification_and_validation
- [31] Design Review, Verification and Validation - Quality Gurus, <https://www.qualitygurus.com/design-review-verification-and-validation/>
- [32] Verification Versus Validation – What’s the Difference? - Climedoo, <http://climedoo.de/en/blog/verification-versus-validation-whats-the-difference/>
- [33] System Validation - SEBoK, https://sebokwiki.org/wiki/System_Validation
- [34] Implementing ISO 15288 V&V Processes using the V&V Studio - The Reuse Company, <https://www.reusecompany.com/wp-content/uploads/2021/02/VV-Studio-Webinar-Jan-2021.pdf>
- [35] Verification (glossary) - SEBoK, [https://sebokwiki.org/wiki/Verification_\(glossary\)](https://sebokwiki.org/wiki/Verification_(glossary))
- [36] “Sapien’s AI Glossary of Data Terms, Definitions & Insights,” Sapien.io, 2025. <https://www.sapien.io/glossary/all>
- [37] A. Pande, “Metamorphic and adversarial strategies for testing AI systems,” Ministry of Testing, Jan 14, 2025. <https://www.ministryoftesting.com/articles/metamorphic-and-adversarial-strategies-for-testing-ai-systems>
- [38] ISO and IEC Make Foundational Standard on Artificial Intelligence Publicly Available, <https://www.holisticai.com/news/iso-iec-22989-foundational-standard-on-ai-open-source>
- [39] The foundational standards for AI | JTC 1, https://jtc1info.org/wp-content/uploads/2022/06/03_08_Paul_Milan_Wei_The-foundational-standards-for-AI-20220525-ww-mp.pdf

- [40] E. Manziuk, O. Barmak, I. Krak, O. Mazurets, and T. Skrypnyk, “Formal Model of Trustworthy Artificial Intelligence Based on Standardization,” CEUR-WS.org, <https://ceur-ws.org/Vol-2853/short18.pdf>
- [41] ISO/IEC TR 24028:2020 - Information technology - Artificial intelligence - OECD.AI, <https://oecd.ai/en/catalogue/tools/isoiec-tr-2402820>
- [42] Exploring the landscape of trustworthy artificial intelligence: Status and challenges, <https://content.iospress.com/articles/intelligent-decision-technologies/1dt240366>
- [43] ISO 42001 Artificial Intelligence Management System - Amazon Web Services (AWS), <https://aws.amazon.com/compliance/iso-42001-faqs/>
- [44] ISO 42001 - AI Management System - BSI, <https://www.bsigroup.com/en-US/products-and-services/standards/iso-42001-ai-management-system/>
- [45] ISO/IEC 22989:2023 Understanding AI Concepts and Definitions Training Course | BSI, <https://www.bsigroup.com/en-ID/training-courses/isoiec-229892023-understanding-ai-concepts-and-definitions-training-course/>
- [46] AI Compliance Audit: Step-by-Step Guide - Dialzara, <https://dialzara.com/blog/ai-compliance-audit-step-by-step-guide/>
- [47] R. Prieto, “Verification and Validation of Project Management Artificial Intelligence Key Points,” Jun. 2020. https://www.researchgate.net/publication/342452507_Verification_and_Validation_of_Project_Management_Artificial_Intelligence_Key_Points
- [48] AI audit checklist (updated 2025) | Complete AI audit procedures | Technical evaluation framework | System reliability guide | Compliance checklist | Lumenalta, <https://lumenalta.com/insights/ai-audit-checklist-updated-2025>
- [49] Y. Wang and S. H. Chung. “Artificial intelligence in safety-critical systems: a systematic review,” Industrial Management & Data Systems,| Emerald Insight, Dec. 2021. <https://www.emerald.com/insight/content/doi/10.1108/imds-07-2021-0419/full/html>
- [50] Trustworthy AI - AI@UCSF - University of California San Francisco, <https://ai.ucsf.edu/trustworthy>
- [51] AI Risks and Trustworthiness - NIST AIRC - National Institute of Standards and Technology, <https://airc.nist.gov/airmf-resources/airmf/3-sec-characteristics/>

- [52] S. A. Seshia, D. Sadigh, and S. Shankar Sastry. Toward Verified Artificial Intelligence. Communications of the ACM, July 2022. <https://cacm.acm.org/research/toward-verified-artificial-intelligence/>
- [53] AI, Opacity, and Personal Autonomy, <https://d-nb.info/1275205275/34>
- [54] Measure - NIST AIRC - National Institute of Standards and Technology, <https://airc.nist.gov/airmf-resources/playbook/measure/>
- [55] Understanding the NIST AI RMF: What It Is and How to Put It Into Practice - Secureframe, <https://secureframe.com/blog/nist-ai-rmfy>
- [56] AI Life Cycle Core Principles - CodeX - Stanford Law School, <https://law.stanford.edu/2023/03/17/ai-life-cycle-core-principles/>
- [57] Ethical and societal implications of algorithms, data, and artificial intelligence: a roadmap for research - Nuffield Foundation, <https://www.nuffieldfoundation.org/sites/default/files/files/Ethical-and-Societal-Impl-ications-of-Data-and-AI-report-Nuffield-Foundat.pdf>
- [58] Messages on “When using AI systems, what are some best practices for ensuring the results you receive are accurate, relevant, and aligned with your original goals?” - ProjectManagement.com, <https://www.projectmanagement.com/discussion-topic/203772/when-using-ai-systems--wha-t-are-some-best-practices-for-ensuring-the-results-you-receive-are-a-correct--relevant--and-aligned-with-your-original-goals-?sort=asc&pageNum=39>
- [59] A Framework for the Verification and Validation of Artificial Intelligence Machine Learning Systems - JagWorks@USA - University of South Alabama, https://jagworks.southalabama.edu/theses_diss/137/y
- [60] Trustworthy AI - The Data Science Institute at Columbia University, <https://datascience.columbia.edu/news/2020/trustworthy-ai/>
- [61] NIST launches ARIA program to assess societal impacts, ensure trustworthy AI systems, <https://industrialcyber.co/ai/nist-launches-aria-program-to-assess-societal-impacts-ensure-trustworthy-ai-systems/>
- [62] User Acceptance Testing (UAT): Definition, Process, and Tools - LambdaTest, <https://www.lambdatest.com/learning-hub/user-acceptance-testing>
- [63] Human-AI Interaction and User Satisfaction: Empirical Evidence from Online Reviews of AI Products - ResearchGate, https://www.researchgate.net/publication/390142284_Human-AI_Interaction_and_User_Satisfaction_Empirical_Evidence_from_Online_Reviews_of_AI_Products/download

- [64] J. Renkhoff, K. Feng, M. Meier-Doernberg, A. Velasquez, and H. H. Song, “A Survey on Verification and Validation, Testing and Evaluations of Neurosymbolic Artificial Intelligence,” *IEEE transactions on artificial intelligence*, pp. 1–15, Jan. 2024, <https://doi.org/10.1109/tai.2024.3351798>.
- [65] A. E. Goodloe, “Assuring Safety-Critical Machine Learning-Enabled Systems: Challenges and Promise,” *Computer*, vol. 56, no. 9, pp. 83–88, Sep. 2023, <https://doi.org/10.1109/mc.2023.3266860>
- [66] A. Woodie, “Top 10 Challenges to GenAI Success,” BigDATAwire, Jan. 22, 2024. <https://www.bigdatawire.com/2024/01/22/top-10-challenges-to-genai-success/>
- [67] C. Bronson, “AI Safety Metrics: How to Ensure Secure and Reliable AI Applications” - Galileo AI, 2025, <https://www.galileo.ai/blog/introduction-to-ai-safety>
- [68] R. Camacho, “A Practical Guide for AI in Safety-Critical Embedded Systems - Parasoft, 2025, <https://www.parasoft.com/blog/ai-in-safety-critical-embedded-systems/>
- [69] Y. Y. Elboher et., al “Formal Verification of Deep Neural Networks for Object Detection,” Arxiv.org, 2023. <https://arxiv.org/html/2407.01295>
- [70] What are non-deterministic AI outputs? - Statsig, 2024, <https://www.statsig.com/perspectives/what-are-non-deterministic-ai-outputs->
- [71] K. Leahy et al., “Grand Challenges in the Verification of Autonomous Systems,” arXiv (Cornell University), Nov. 2024, [https://doi.org/10.48550/arxiv.2411.14155.](https://doi.org/10.48550/arxiv.2411.14155)
- [72] Symbolic AI vs. Deep Learning: Key Differences and Their Roles in AI Development, <https://smythos.com/ai-agents/agent-architectures/symbolic-ai-vs-deep-learning/>
- [73] Symbolic AI vs. Machine Learning: A Comprehensive Guide - SmythOS, <https://smythos.com/ai-agents/ai-tutorials/symbolic-ai-vs-machine-learning/>
- [74] O. Vermesan, V. Piuri, F. Scotti, A. Genovese, R. D. Labati, and P. Coscia, “Explainability and Interpretability Concepts for Edge AI Systems,” River Publishers eBooks, pp. 197–227, Feb. 2024, [https://doi.org/10.1201/9781003478713-9.](https://doi.org/10.1201/9781003478713-9)
- [75] What Is Explainable AI (XAI)? Palo Alto Networks, <https://www.paloaltonetworks.com/cyberpedia/explainable-ai>
- [76] Deep Learning’s Challenges and Neurosymbolic AI’s Solutions - AskUI, 2024. <https://www.askui.com/blog-posts/deep-learnings-challenges-and-neurosymbolic-ais-solutions>

- [77] V. Musanga, S. Viriri, and C. Chibaya, “A Framework for Integrating Deep Learning and Symbolic AI Towards an Explainable Hybrid Model for the Detection of COVID-19 Using Computerized Tomography Scans,” *Information*, vol. 16, no. 3, p. 208, Mar. 2025, <https://doi.org/10.3390/info16030208>.
- [78] X. Zhang, “Apex.OS: Breaking Barriers in Autonomous Verification & Validation”, 2024, <https://www.apex.ai/post/apex-os-breaking-barriers-in-autonomous-verification-validation>
- [79] J. Szarmach, NIST: Reducing Risks Posed by Synthetic Content An Overview of Technical Approaches to Digital Content Transparency, 2025, <https://www.aigl.blog/nist-reducing-risks-posed-by-synthetic-content-an-overview-of-technical-approaches-to-digital-content-transparency/>
- [80] NIST Trustworthy and Responsible AI - NIST AI 100-4. Reducing Risks Posed by Synthetic Content. An Overview of Technical Approaches to Digital Content Transparency. <https://nvlpubs.nist.gov/nistpubs/ai/NIST.AI.100-4.pdf>
- [81] NIST. The United States Artificial Intelligence Safety Institute: Vision, Mission, and Strategic Goals, 2024, <https://www.nist.gov/document/ai-si-strategic-vision-document>
- [82] S. Dahaweer, “How AI is going to revolutionize safety in critical applications”, Alithya, 2023, <https://www.alithya.com/en/insights/blog-post/how-ai-going-revolutionize-safety-critical-applications>
- [83] The Top 10 Unsolved Challenges in AI: A 2024 Retrospective, gekko, 2024, <https://gpt.gekko.de/unsolved-challenges-in-ai-2024/>
- [84] Risks from AI - An Overview of Catastrophic AI Risks, CAIS - Center for AI Safety, <https://www.safe.ai/ai-risk>
- [85] R. Lubecki, “Verifying and Validating AI/ML” - UpCity, 2021, <https://upcity.com/experts/verifying-and-validating-ai-ml/>
- [86] K. Lekadir et al., “FUTURE-AI: International consensus guideline for trustworthy and deployable artificial intelligence in healthcare,” *BMJ*, vol. 388, p. e081554, Feb. 2025, <https://doi.org/10.1136/bmj-2024-081554>.
- [87] Traditional AI vs. Generative AI: What’s the Difference? - College of Education, Illinois, 2024, <https://education.illinois.edu/about/news-events/news/article/2024/11/11/what-is-generative-ai-vs-ai>
- [88] N. Ahsan, “Why Enterprises Are Adopting Domain-Specific AI Agents”, Vidizmo, 2025, <https://vidizmo.ai/blog/why-domain-specific-ai-agents-are-key-to-business-success>

- [89] Future of AI Research – AAAI – Association for the Advancement of Artificial Intelligence 2025. Available: <https://aaai.org/wp-content/uploads/2025/03/AAAI-2025-PresPanel-Report-FINAL.pdf>
- [90] Methods For Verifying AI Reliability to Ensure AI Safety, IAI – Institute for AI Transformation, August 2024. Available: <https://www.leadersinai summit.com/insights/methods-for-verifying-ai-reliability-to-ensure-ai-safety>
- [91] IEEE 1012-2016. IEEE Standard for System, Software, and Hardware Verification and Validation. <https://webstore.ansi.org/standards/ieee/ieee10122016?source=blog>
- [92] O. Vermesan, K. De Bosschere, T. Vardanega, S. Azaiez, M. Duranton, R. Badia, and D. Lezzi (Eds.). “Distributed Computing and Swarm Intelligence - Developing a Vision for Transatlantic Collaboration,” Zenodo, Feb. 2025, <https://doi.org/10.5281/zenodo.14940197>

2

Pioneering the Hybridization of Federated Learning in Human Activity Recognition

**Alfonso Esposito¹, Yasamin Moghbelan¹, Ivan Zyrianoff¹,
Leonardo Ciabattini¹, Federico Montori^{1,2}, and Marco Di Felice^{1,2}**

¹University of Bologna, Italy

²University of Bologna, Advanced Research Center for Electronic Systems, Italy

Abstract

The Internet of Things (IoT) nowadays greatly benefits from Artificial Intelligence (AI) algorithms implemented in the edge, because of their efficiency and the reduction of costs that they imply. The advent of Federated Learning (FL) has made possible the combination of the advantages of edge-AI, among which the privacy of users, as their data is not shared with the cloud, with the collective intelligence. However, FL is known to have worse performance compared to its centralized counterpart, which may not be tolerable in certain cases. In this paper, we propose a hybrid framework for FL, imagining a number of clients that are willing to share part of their data. We envisioned two types of Hybridization: vertical and horizontal. The goal of this paper is to assess whether a small hybridization can bring advantages to the overall performance of the whole FL procedure in terms of classification accuracy.

Keywords: Internet of Things (IoT), deep learning (DL), federated learning (FL), human activity recognition (HAR), performance evaluation.

2.1 Introduction and Background

The growth of the Internet of Things (IoT) has enabled novel monitoring systems capable of gathering data from heterogeneous and pervasive devices, supporting smart city, healthcare and industrial domains. This data

is processed using advanced Deep Learning (DL) techniques for system state forecasting, contributing to the integrated paradigm of the Artificial Intelligence of Things (AIoT) [1]. However, it is well known that advanced DL techniques, particularly those used in computer vision applications, require large datasets with adequate number of instances for each system state or class to be predicted. In many cases, building reliable DL models requires aggregating data from multiple heterogeneous users or organizations performing decentralized data collection for a common task. This is the case, for instance, of most of Human Activity Recognition (HAR) applications that utilize DL models trained from wearable or camera-based IoT data gathered from multiple volunteers or via crowdsensing techniques [2].

State-of-the-art AIoT systems often rely on cloud-based centralized architectures, which facilitate the aggregation of IoT data from diverse clients. While these solutions enable easy deployment and scalable computational and storage resources, they also raise significant privacy concerns due to the inherent risks associated with data sharing. Federated Learning (FL), first proposed in 2018 [3], has emerged as a privacy-preserving alternative to centralized machine learning, enabling cooperative training in a distributed environment. In a typical FL setup, multiple clients independently train models on their private datasets and only share the trained weights with a central server, which aggregates these weights and sends the updated model back to the clients. This approach ensures that no raw data is exchanged among clients, although concerns about the trustworthiness of the centralized server remain [4]. Moreover, variations in data quality and quantity across clients can create challenges in ensuring fair evaluation of each client's contribution and achieving a high-quality global model [5]. To address the issue of non-i.i.d. (non-independent and identically distributed) data across clients, various solutions have been proposed, such as clustering clients with similar data distributions or adopting weight-based model aggregation techniques [6].

In this paper, we aim to bridge the gap between centralized and Federated Learning (FL) methodologies, in order to address distributed IoT use cases where privacy requirements vary from client to client. For instance, some clients may be willing to share raw data, while others may not, due to differences in client nature (e.g., public vs. private organizations), varying perceptions of privacy—often shaped by social factors [7]—or monetization strategies, where some clients are incentivized to sell their data. We refer to this scenario as *hybrid FL* and describe this *hybridization* approach for data gathering and model building at the server side, which offers a novel interpretation compared to other studies [4]. Specifically, this paper

explores two types of FL hybridization: *vertical* and *horizontal*. In the vertical hybridization, a portion of clients shares raw data with the centralized server, while others only share deep learning model coefficients generated from local training. In the horizontal case, all clients share a variable portion of their raw data (e.g., corresponding to the amount for which they have been compensated) in addition to the DL model coefficients trained on the complementary data. We evaluate the performance of both strategies using two benchmark datasets (related to Human Activity Recognition and image classification) and analyze the impact of different levels of hybridization compared to pure centralized and FL-based approaches. Our results demonstrate that hybridization can be a powerful tool for improving the accuracy of federated systems, even when applied slightly.

To summarize, the key contributions of our paper are the following:

- Introduction of hybrid FL as a new strategy for privacy-adaptive learning in IoT scenarios.
- Proposal of two versions of hybrid FL, respectively horizontal and vertical, based on distinct methods of integrating raw data and deep learning weights at the server.
- Evaluation of proposed strategies across varying degrees of hybridization, using two widely adopted benchmark datasets in the DL community.

The rest of the paper is structured as follows. Section 2 introduces the revised FL architecture supporting the vertical and horizontal hybridization. Section 3 describes the evaluation methodology, datasets and metrics. Section 4 presents some evaluation results. Section 5 concludes the work and discusses future research steps.

2.2 Hybrid FL Architecture

We consider the scenario depicted in Figure 2.1, composed of two main actors: N distributed clients and the server. Each client c_i possesses its own dataset D_i gathered through its local sensors, and a DL network topology, denoted as m in the following. In a classical FL application, each client c_i performs local training rounds of model m on D_i and then shares the list of weights W_i with the server: the latter is responsible for aggregating the weights, for instance through the FedAvg [3] algorithm, and sending the updated values back to the clients.

In the proposed hybrid FL architecture, the server performs additional storage and computational tasks, basically behaving as a special client device.

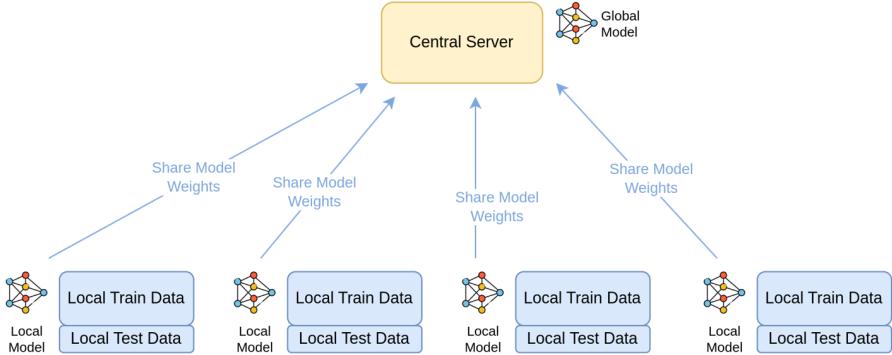


Figure 2.1 Typical Federated Learning Architecture

Indeed, it possesses a local dataset D_s that is built from clients' datasets or portions of them. More formally, we have that $D_s = \bigcup_{i=0}^n D_i^*$ where $n \leq N$ and $D_i^* \subseteq D_i$. The server trains the m model on D_s getting a local version of weights W_S that is later averaged with the weights W_i received by clients at each round. We denoted this process as **FL hybridization**, distinguishing between two modes for creating the local dataset D_s :

- **Vertical Hybrid FL.** In such case, we have that: $n < N$ and $D_i^* = D_i$. In other words, only a subset of the client nodes shares its own raw data with the server. This setup may model two different use-cases: (i) only n clients have been compensated with monetary rewards to share their data and/or (ii) n clients do not consider their local datasets privacy-sensitive. We indicate with $r_v = \frac{n}{N}$ the rate of vertical hybridization. Clearly, for $n = N$ and $r_v = 1$, the system is equivalent to centralized learning. Vice versa, for $n = 0$ and $r_v = 0$, the system is equivalent to a pure FL approach. We investigate the impact of varying r_v configurations in the overall DL performance in Section 4. We show in Figure 2.2 an overview on the conceptual architecture of Vertical Hybrid FL.
- **Horizontal Hybrid FL.** In such case, we have that: $n = N$ and $D_i^* \subset D_i$. In other words, all clients share only a portion of their local datasets with the server. This setup may model a practical use-case where each client shares with the server an amount of raw data proportionally to the monetary reward it received. We indicate with $r_h = \text{mean}_{0 \leq i < N} \left(\frac{|D_i^*|}{|D_i|} \right)$ the rate of horizontal hybridization. Clearly, for $D_i^* = D_i$, the system becomes a centralized learning. Vice versa, if $D_i^* = \emptyset \forall c_i$ the system is

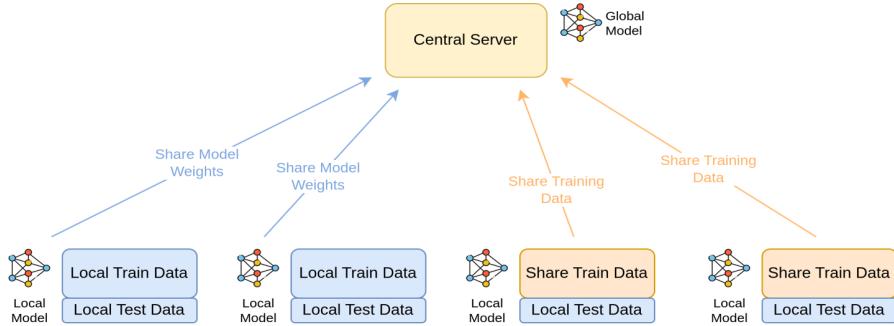


Figure 2.2 Vertical Hybrid Federated Learning Architecture

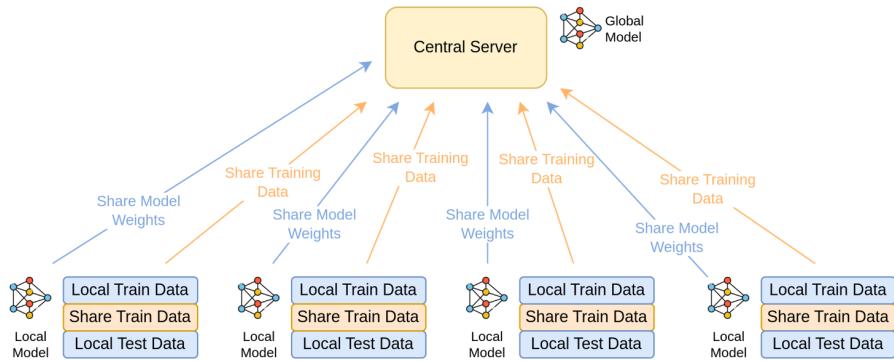


Figure 2.3 Horizontal Hybrid Federated Learning Architecture

equivalent to a pure FL approach. We investigate the impact of varying r_h configurations in the overall DL performance in Section 4. We show in Figure 2.3 an overview on the conceptual architecture of Vertical Hybrid FL.

We further highlight that the hybridization rates determine the amount of required privacy preservation. The maximum value (1) corresponds to when all clients share their local data with the server. Vice versa, the minimum value (0) forbids any transmission of raw data outside the clients' devices.

2.3 Evaluation Methodology and Metrics

In this section, we describe the experiments that we performed to investigate the performances of the two hybrid approaches explained in the previous

section. First, we will describe the methodology used to test the effectiveness of HFL across two datasets: FEMNIST and UCI HAR.

The widely recognized University of California Irvine (UCI) HAR dataset [8] was created using data from smartphone accelerometer and gyroscope sensors, which were used to classify six types of human activities. It was gathered from 30 volunteers, each carrying a smartphone while performing six distinct activities: walking, walking upstairs, walking downstairs, sitting, standing, and lying down. The dataset consists of time-series data captured across the three axes of both sensors, along with the corresponding activity labels. It has been extensively used in research for developing and assessing HAR models using machine learning techniques. Each record in the dataset contains a vector of 561 features, derived from both time and frequency domain calculations.

The FEMNIST dataset [9] is an adaptation of the extended version of the MNIST dataset that has been modified to be suitable for FL tasks. The MNIST dataset contains 28 by 28 pixels images of handwritten digits and characters (62 classes in total), and the goal of the DL model is to guess the actual character represented. The FEMNIST dataset groups the elements on top of the user that actually performed the handwriting, producing a number of sub-datasets each of them with a different style of writing. The number of users of the FEMNIST dataset is 3500, however, for the purpose of our experiment, we considered only 30 users, in order to make experiments comparable between the two datasets.

For UCI HAR we employed, as a base local model, a simple feed-forward neural network, while for FEMNIST we adopted a convolutional network. Each of the local sub-datasets is split into training and test set using a stratified split with a 70%-30% ratio.

We performed federated classification experiments by employing 6 epochs and 20 rounds of federation, recording the accuracy score at the end of the last round.

We tested both vertical and horizontal hybridization, by setting alternately r_v and r_h to values spanning from 0% to 100% with a 10% step.

The experiments were implemented in Python using the Flower framework (<https://flower.readthedocs.io/en/latest/>). Since Flower does not support the implementation of hybridization, we adopted the two following methods to simulate the two hybridization methods (as Figure 2.4 suggests):

- Vertical Hybridization was simulated by aggregating all the clients that share their whole dataset instead of the weights into a single client.

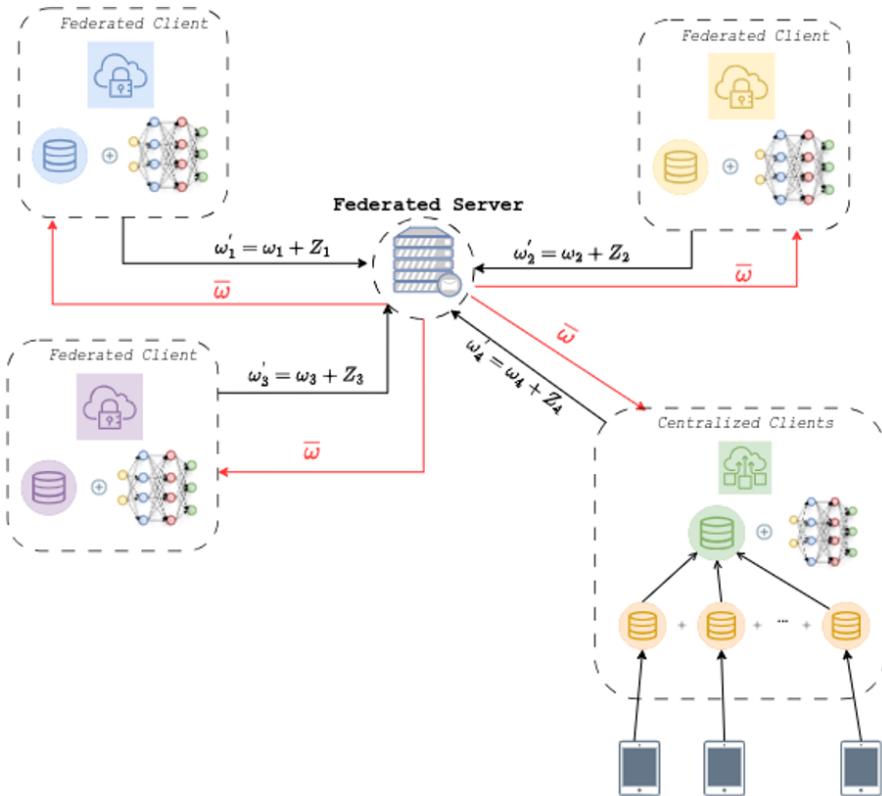


Figure 2.4 Implementation of the Vertical Hybridization in Flower

- Horizontal Hybridization was simulated by extracting from each client the portion of the training dataset that they aim to share and assigning it to a new “sink” client.

Each experiment was then repeated 20 times, by randomizing the clients or the dataset portions to be shared. This ensures scientific rigor and smooths out certain corner situations that may arise.

2.4 Evaluation Results

This section presents the outcome of the experiments presented in the previous section. The results are aimed at evaluating the HFL approaches on the datasets. We investigate how different levels of data sharing in vertical

and horizontal hybridization affect model performance. We first discuss the results for the UCI HAR dataset, followed by FEMNIST, to uncover any dataset-specific trends and performance differences. A general overview of the results shows, as expected, an overall improvement in model performance as the degree of hybridization increases, which is notable for both hybrid approaches. The UCI HAR dataset performed best. Even with the relatively small model size, it consistently achieved excellent results, maintaining accuracy above 90% and reaching almost 95% in experiments with higher levels of hybridisation. This is shown in Figures 2.5 and 2.6, where the increase in model performance as the level of hybridisation increases is evident, with an increase of 2 percentage points already at a low level of horizontal hybridization (20%).

The FEMNIST is the dataset where the performance improvements from sharing data is most noticeable. As we can observe in the Figures 2.7 and 2.8, the sharing of a small number of data points could lead to a significant improvement in performance. Specifically, sharing 10% of the data led to an improvement of approximately 5% in accuracy, while sharing 20% led to an additional improvement of approximately 2/3, resulting in a final performance of 88%. This is comparable to the 90% accuracy obtained through centralized training. Beyond a data sharing rate of 30%, the improvements

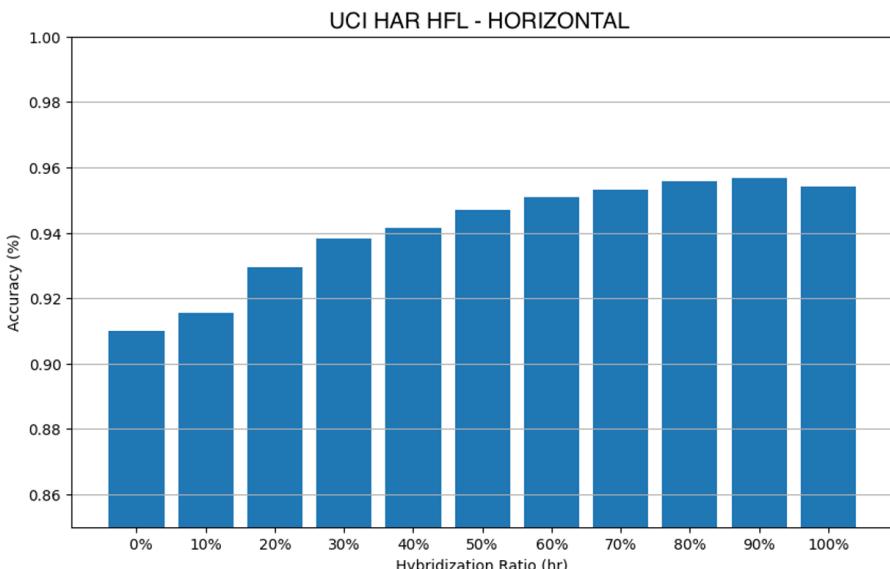


Figure 2.5 Horizontal Hybridization Results for UCI HAR

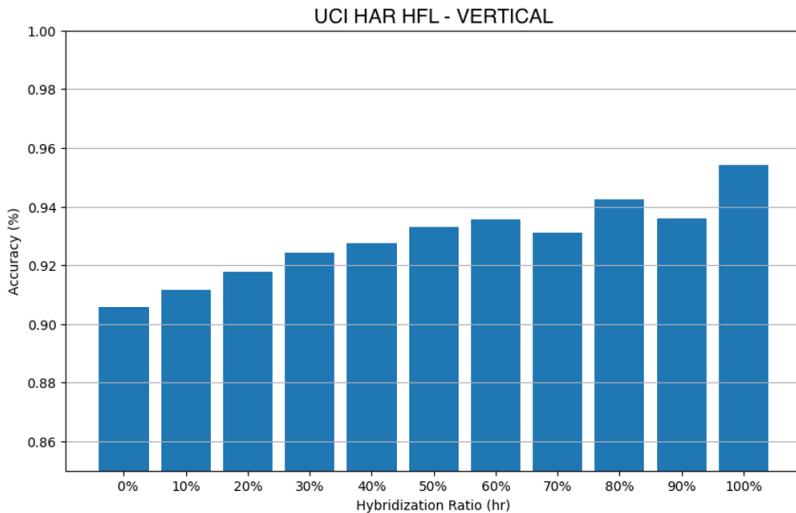


Figure 2.6 Vertical Hybridization Results for UCI HAR

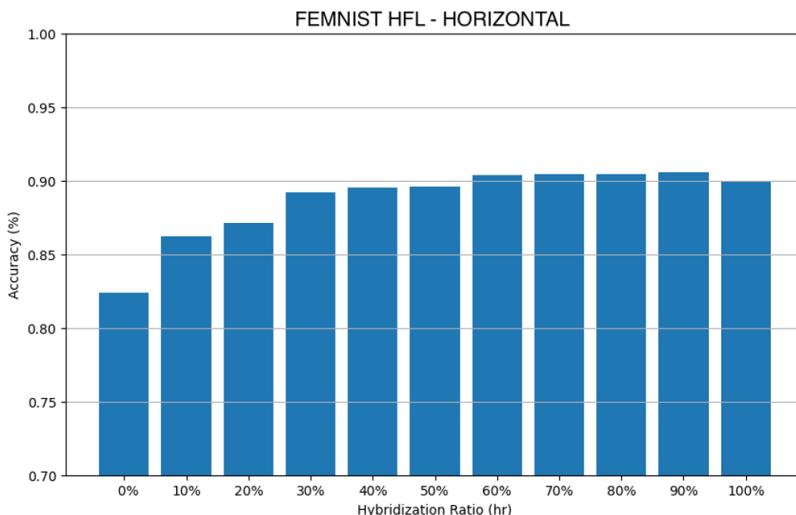


Figure 2.7 Horizontal Hybridization Results for FEMNIST

obtained are increasingly marginal, with a maximum of 1%. This suggests that a data sharing rate of 20% represents an optimal balance between performance and data sharing.

The results demonstrate that the sharing of a portion of the dataset has a significant positive impact on the model's performance. This effect was

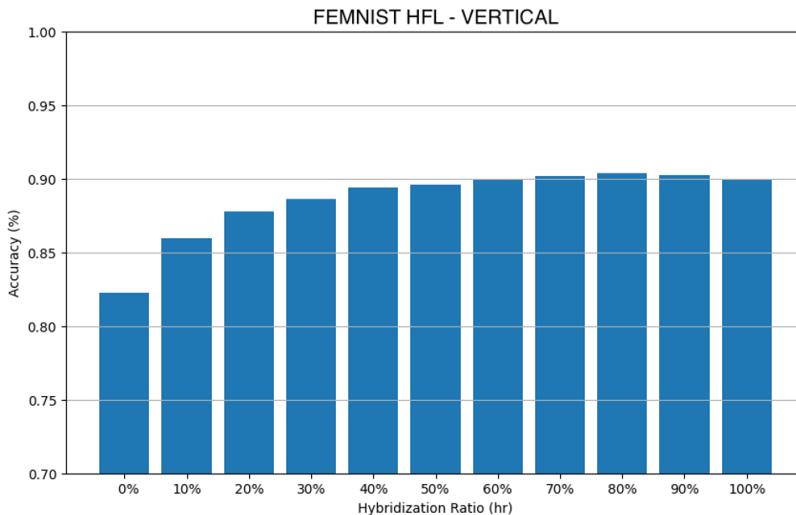


Figure 2.8 Vertical Hybridization Results for FEMNIST

observed across two distinct datasets, UCI HAR and FEMNIST, indicating that the improvement is not specific to any problem or model. The performance enhancement was especially evident in the case of the FEMNIST dataset.

2.5 Conclusion and Future Works

In this paper we examined the effects of hybridization for Federated Learning scenarios. We specifically directed our research towards Human Activity Recognition, imagining scenarios in which certain clients would be willing to share (part of) their data with the central server for a reward, penalizing their privacy to an extent. Results showed that a minimal amount of hybridization does provide an increase in the performance. The extent to which the privacy of the users is compromised by this is a future work. We aim to study how to select carefully data in a way in which the privacy is minimally affected, as well as to blend the two hybridization techniques, to select the best configuration.

Acknowledgements

This research is funded by the “Progetto Casa delle Tecnologie Emergenti” - Comune di Bologna - PSC MISE 2014-2020.

References

- [1] K. S. Awaisi, Q. Ye and S. Sampalli, “A Survey of Industrial AIoT: Opportunities, Challenges, and Directions,” in *IEEE Access*, vol. 12, pp. 96946-96996, 2024. <https://doi.org/10.1109/ACCESS.2024.3426279>
- [2] S. Ankalaki, “Simple to Complex, Single to Concurrent Sensor-Based Human Activity Recognition: Perception and Open Challenges,” in *IEEE Access*, vol. 12, pp. 93450-93486, 2024. <https://doi.org/10.1109/ACCESS.2024.3422831>
- [3] H. Brendan McMahan, E. Moore, D. Ramage, S. Hampson, B. Agüera y Arcas, “Communication-Efficient Learning of Deep Networks from Decentralized Data”, arXiv, <https://arxiv.org/abs/1602.05629>
- [4] J. Cui, H. Zhu, H. Deng, Z. Chen, and D. Liu, “Fearh: Federated machine learning with anonymous random hybridization on electronic medical records,” *Journal of Biomedical Informatics*, vol. 117, p. 103735, 2021. <https://doi.org/10.1016/j.jbi.2021.103735>
- [5] M. Moshawrab, M. Adda, A. Bouzouane, H. Ibrahim, and A. Raad, “Reviewing federated learning aggregation algorithms; strategies, contributions, limitations and future perspectives,” *Electronics*, vol. 12, no. 10, p. 2287, May 2023. <https://doi.org/10.3390/electronics12102287>
- [6] H. Lee and D. Seo, “FedLC: Optimizing Federated Learning in Non-IID Data via Label-Wise Clustering,” in *IEEE Access*, vol. 11, pp. 42082-42095, 2023. <https://doi.org/10.1109/ACCESS.2023.3271517>
- [7] D. Ibdah, N. Lachtar, S. M. Raparthi and A. Bacha, “Why Should I Read the Privacy Policy, I Just Need the Service”: A Study on Attitudes and Perceptions Toward Privacy Policies,” in *IEEE Access*, vol. 9, pp. 166465-166487, 2021. <https://doi.org/10.1109/ACCESS.2021.3130086>
- [8] Anguita, Davide, Alessandro Ghio, Luca Oneto, Xavier Parra, and Jorge Luis Reyes-Ortiz. “A public domain dataset for human activity recognition using smartphones.” In *Esann*, vol. 3, p. 3. 2013. <https://www.esann.org/sites/default/files/proceedings/legacy/es2013-84.pdf>
- [9] Caldas, Sebastian, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H. Brendan McMahan, Virginia Smith, and Ameet Talwalkar. “Leaf: A benchmark for federated settings.” *arXiv preprint arXiv:1812.01097*. <https://arxiv.org/abs/1812.01097>