



PhD in Information Technology and Electrical Engineering
Università degli Studi di Napoli Federico II

PhD Student: Antonio Emmanuele

Cycle: XXXIX

Training and Research Activities Report

Academic year: 2024-25 - PhD Year: Second

Student Signature: *Antonio Emmanuele*

Tutor: prof. Mario Barbareschi

Tutor Signature: *Mario Barbareschi*

Co-Tutor: -

Date: 31/10/2025

Training and Research Activities Report

PhD in Information Technology and Electrical Engineering

Cycle: XXXIX

Author: Antonio Emmanuele

1. Information:

- **PhD student:** Antonio Emmanuele **PhD Cycle:** XXXIX
- **DR number:** DR997187
- **Date of birth:** 04/12/1998
- **Master Science degree:** Computer Engineering **University:** Federico II of Naples
- **Scholarship type:** UNINA
- **Tutor:** Mario Barbareschi
- **Co-tutor:** -
- **Period abroad:** Institut des Nanotechnologies de Lyon (INL) of the University École Centrale de Lyon, spent period: 07/10/2025 to 31/11/2025 (2nd year), remaining period (1/11/2025 – 6/04/2026) (3nd year).

2. Study and training activities:

Activity	Type ¹	Hours	Credits	Dates	Organizer	Certificate ²
Safety Critical Systems for Railway Traffic Management	Course	17	4	3/10/2025 – 30/10/2025	ITEE, DIETI	Y
1 st IWES PhD School on Embedded Systems	Doctoral School	24	4.8	15-17 September 2025	CINI ESSM Laboratory and University of Modena and Reggio Emilia Modena (Italy)	Y
20th TAROT Summer School on Software Testing, Verification & Validation	Doctoral School	24	4.8	June 30 - July 4, 2025	University Federico II of Naples	Y

1) Courses, Seminar, Doctoral School, Research, Tutorship

2) Choose: Y or N

Training and Research Activities Report

PhD in Information Technology and Electrical Engineering

Cycle: XXXIX

Author: Antonio Emmanuele

2.1. Study and training activities - credits earned

	Courses	Seminars	Research	Tutorship	Total
Bimonth 1			9		9
Bimonth 2			9		9
Bimonth 3			10		10
Bimonth 4			10		10
Bimonth 5		4.8	2	3.2	10
Bimonth 6	4	4.8	5		13.8
Total	4	9.6	45	3.2	61.8
Expected	30 - 70	10 - 30	80 - 140	0 - 4.8	

3. Research activity:

Decision Tree Based Models, Approximation and Acceleration

The increasing adoption of the Internet of Things, along with the significant amount of data it generates, is fostering the widespread use of the edge-computing paradigm. This paradigm envisions that computation is performed directly on nodes located at the edge of the computing infrastructure. By avoiding data transmission to remote cloud servers, edge computing provides substantial advantages in terms of (1) latency, (2) network bandwidth, and (3) end-user privacy. In this context, due to its widespread use, Machine Learning has become one of the most common tasks offloaded to edge nodes. However, edge devices are highly heterogeneous, including, for instance:

- Field Programmable Gate Arrays (FPGAs):** These are limited in capability when a Machine Learning model is directly translated into a hardware implementation, as the model size directly impacts the hardware resources required by the accelerator.
- Microcontrollers:** These integrate a classical microprocessor architecture along with several peripherals within the same chip. Such nodes lack the memory and computational capabilities of fully-fledged machines. Moreover, they are typically battery-powered, requiring task execution with minimal energy consumption.
Notably, as edge computing becomes more pervasive, these nodes are starting to include richer memory hierarchies—such as caches and Tightly Coupled Memories (TCMs)—as well as advanced processor extensions like Single Instruction Multiple Data (SIMD). Despite this trend, existing algorithms for edge tasks, such as Machine Learning models, still require adaptation. Indeed, the computational model of these embedded edge nodes differs from that of fully-fledged machines for which compute-intensive tasks are typically designed. For example, high-end microprocessors generally do not include Tightly Coupled Memory, which can be leveraged to design algorithms that avoid cache-related uncertainty. This is because, unlike caches, TCMs are memories explicitly addressable by the executing code.

Training and Research Activities Report

PhD in Information Technology and Electrical Engineering

Cycle: XXXIX

Author: Antonio Emmanuele

In this scenario, Deep Learning models are no longer considered a universal solution. On one hand, their execution demands considerable computational resources, increasing the burden on edge nodes. On the other hand, it has already been shown that several tasks—such as time series forecasting and point-wise anomaly detection—can be efficiently addressed using Decision Tree-based models. These models, usually ensembles of Decision Trees, consist of binary trees where each leaf represents a specific inference outcome. Their inference process corresponds to a simple tree traversal, without the need for complex matrix multiplications as in Convolutional Neural Networks (CNNs).

Despite these advantages, such models cannot be seamlessly deployed at the edge. Specifically:

1. **FPGA implementations:** Translating these models into dedicated hardware accelerators leads to scalability issues. As training datasets and task complexity increase, model size also grows, resulting in unrealistic hardware requirements. The corresponding FPGA area occupancy may exceed the capacity typically available on Commercial Off-The-Shelf (COTS) devices.
2. **Reliability of FPGA Accelerators:** The adoption of such models at the edge, especially in critical domains such as automotive, requires building trust in their hardware implementations. In this context, it is essential to understand how these circuits fail and to estimate the probability that a transient or permanent fault compromises their accuracy. Unfortunately, the scientific literature still lacks comprehensive studies addressing the reliability of Machine Learning models deployed at the edge.
3. **Microcontroller implementations:** Despite the growing computational capabilities of microcontrollers, these models cannot fully exploit the available hardware resources. Decision Trees are irregular data structures that, as established in the High-Performance Computing community, cannot efficiently leverage caching or SIMD extensions.

In light of these challenges, during my second year I investigated the following topics:

1. **Approximation of FPGA acceleration for Machine Learning models.**
Approximate Computing (AxC) is a design paradigm based on the observation that many computing applications can still deliver accurate results without performing every computation with maximum precision. Based on this principle, AxC can be applied to Machine Learning models, which typically exhibit a high degree of redundant computation that can be reduced to obtain approximate yet effective models.
2. **Resiliency Analysis of Decision Tree Ensembles**
As for model-specific FPGA accelerators, it remains unclear how these circuits behave in the presence of faults. A widely used approach to assess reliability is fault injection, which consists of introducing a transient or permanent fault into a system and evaluating its impact on the final outcome. However, extensively characterizing Decision Tree accelerators through fault injection is challenging due to the large number of potential fault sites. To

address this, I employed the Statistical Fault Injection technique, a well-established method in the circuit reliability community that considers only a statistically significant subset of possible fault locations. Using this approach, I investigated the reliability of Decision Tree accelerators by analyzing their resilience (i) when faults are injected into different hardware components, and (ii) when approximation is applied.

3. **SIMD execution of Decision Tree-based models.**

Due to the limited availability of SIMD extensions in resource-constrained edge nodes, SIMD-based visiting procedures have rarely been considered for Decision Tree inference. In view of recent architectural trends, I proposed a tree-based inference algorithm tailored for SIMD execution. This approach leverages Tightly Coupled Memories, which, unlike caches that are difficult to exploit for irregular data structures, exhibit deterministic behavior and therefore enable full utilization of SIMD extensions.

Hardware Security Based On Physical Unclonable Functions.

Due to the widespread diffusion of edge nodes—typically resource-constrained devices—it is crucial to design and implement strong security mechanisms among them. The main security challenges stem from the intrinsic characteristics of edge nodes, which are limited in memory, computational power, and energy budget. Moreover, since these devices are often deployed directly in the field, they are easily exposed to physical tampering.

In summary, security solutions for edge nodes must satisfy the following constraints:

1. **Minimal reliance on pre-stored cryptographic material**, typically stored in costly secure memories that are often unavailable in resource-limited devices.
2. **Use of computationally lightweight cryptographic algorithms**, to minimize the processing overhead of security tasks.
3. **Linear scalability** in terms of network bandwidth and computational requirements, as edge environments usually involve large groups of collaborating IoT nodes.

The main security tasks required at the edge include:

1. **Efficient key sharing**, enabling the exchange of cryptographic material with minimal computational and communication overhead in both end-to-end and group settings.
2. **Remote attestation**, allowing an edge node to prove to an external verifier that its internal state (e.g., the executing code) has not been tampered with. This requires the device to measure its internal state and produce an **Attestation Report**, which must be tamper-evident, enabling a remote verifier to validate its integrity and detect any modification by an attacker.

In parallel, it is essential that these security mechanisms operate correctly even in **multi-user environments**, where multiple heterogeneous applications—possibly from different

developers—are deployed on the same edge node. This approach enhances flexibility and reduces the hardware overhead of security mechanisms by sharing resources.

With this objective, **Physical Unclonable Functions (PUFs)** represent a lightweight and effective solution to address these challenges. PUFs are silicon circuits capable of generating cryptographic keys by exploiting manufacturing process variations intrinsic to each chip. This randomness provides strong tamper resistance, allowing a verifier to detect any physical alteration of the circuit.

During my second year, I investigated the following topics:

1. **PUF-based Group Remote Attestation:** Study of how PUFs can be used to attest the integrity of groups of edge nodes with minimal network overhead. The main challenge lies in designing a group attestation mechanism that enables (i) efficient aggregation of attestation reports and (ii) rapid identification of compromised nodes when integrity verification fails, with minimal message exchanges.
2. **PUFs in multi-user settings:** Traditional PUF-based protocols assume exclusive ownership of the PUF by a single user. My research focused on the **virtualization of PUFs (vPUFs)**, a security abstraction that provides multiple virtual PUF instances to different users while maintaining strong security isolation and preventing interference among them.

4. Research products:

1- A Decentralized PUF-Based Scheme for Remote Attestation

Status: Accepted and Published

DOI: https://doi.org/10.1007/978-3-032-00644-8_10

Venue: ARES Conference, STAM Workshop

Authors: Mario Barbareschi, Antonio Emmanuele, Daniele Lombardi

Indexes: Scopus

Abstract: With the continuous growth of edge computing, numerous edge devices collaborate in decentralized groups to provide flexible and reliable services while distributing computational workloads. However, the presence of numerous devices requires the adaptation of existing end-to-end security-procedures designed for resource-constrained edge-nodes. Among them, Remote Attestation allows a remote server to directly evaluate the trustworthiness of a node, by remotely verifying that no-malicious code is being executed on

Training and Research Activities Report

PhD in Information Technology and Electrical Engineering

Cycle: XXXIX

Author: Antonio Emmanuele

the target. Unfortunately, attesting groups of devices remains an open challenge, as it must minimize message exchanges to reduce network bandwidth consumption while ensuring lightweight computation on both the server and device sides. For these reasons, we propose DHERAP, a PUF-based remote attestation protocol targeted for groups of decentralized edge-nodes. We design DHERAP to reduce message overhead by enabling different nodes to attest local groups of nodes, removing the need for remote communication to a server.

2- Bridging Efficient and Explainable Traffic Flow Prediction on the Edge

Status: Accepted and Published

DOI: https://doi.org/10.1007/978-3-031-87778-0_34

Venue: International Conference on Advanced Information Networking and Applications

Indexes: Scopus

Authors: Mario Barbareschi, Antonio Emmanuele, Nicola Mazzocca, Franca Rocco Di Torrepadula

Abstract: Digital transformation is significantly reshaping the way transportation systems operate, bringing both new opportunities and challenges. Indeed, Intelligent Transportation Systems (ITSs) are facing accurate traffic flow prediction as a fundamental means for congestion mitigation, route planning, and dynamic traffic signal control, by leveraging machine learning techniques. While deep learning models, like Recurrent Neural Networks, have proved strong predictive performance, they require substantial computational resources and lack interpretability, hindering back their adoption on low-end edge devices and failing to provide required robustness of real-world applications, as required by regulatory standards. Bearing in mind the above, in this paper, we devise an approach that combines local execution onto edge devices of traffic flow prediction by using XGBoost, a tree-ensemble machine learning model, with SHapley Additive exPlanations (SHAP) technique, which guarantees explainability. We detail the whole flow and prove its effectiveness by running it against traffic dataset, showing that we are able to achieve competitive accuracy and local-interpretation of predictions.

3- PUF-Based Secure Key Management for Continuum Computing

Status: Accepted and Published

Training and Research Activities Report

PhD in Information Technology and Electrical Engineering

Cycle: XXXIX

Author: Antonio Emmanuele

DOI: https://doi.org/10.1007/978-3-031-87778-0_38

Venue: International Conference on Advanced Information Networking and Applications

Authors: Mario Barbareschi, Valentina Casola, Antonio Emmanuele, Daniele Lombardi

Indexes: Scopus

Abstract: In the domain of continuum computing, managing security keys efficiently and securely is paramount due to the diverse range of devices involved, from resource-constrained edge devices to powerful cloud servers. This paper introduces an extended version of the Group-Key PHEMAP protocol, initially designed for IoT applications, tailored to the continuum computing environment. The protocol exploits Physical Unclonable Functions (PUFs) for lightweight and secure key generation, maintaining low computational and communication overheads. It supports dynamic group membership, allowing devices to join and leave the group seamlessly without the need for additional cryptographic keys. Performance evaluations are conducted using an advanced network simulator and prototype implementations on devices representative of continuum computing, confirming the protocol feasibility and robustness.

4- Harnessing Explainable AI in Railway: A Decision Tree-Based Approach

Status: Accepted and Published

DOI: <https://doi.org/10.1109/EDCC-C66476.2025.00043>

Venue: European Dependable Computing Conference (EDCC), AI4RAILS workshop

Authors: Mario Barbareschi, Antonio Emmanuele, Nicola Mazzocca, Franca Rocco Di Torrepadula

Indexes: Scopus

Abstract: In recent years, Artificial Intelligence has gained significant popularity for solving various tasks, including service optimization, system monitoring, and industrial control. Despite its success, adoption in critical systems, such as the railway domain, remains limited. This is primarily due to the high stakes in these systems, where failures can lead to damage to critical infrastructure and risks to human lives. As a result, software in these domains must be deterministic, ensuring that all behaviors can be statically verified. Machine Learning models, due to their complexity, are often perceived as black-box systems and exhibit seemingly nondeterministic behavior, making their integration into such infrastructure challenging. To

Training and Research Activities Report

PhD in Information Technology and Electrical Engineering

Cycle: XXXIX

Author: Antonio Emmanuele

address this issue, one potential solution is the use of eXplainable Artificial Intelligence (XAI) techniques, which enable the construction of human-interpretable explanations for model predictions. In this paper, we propose a time-series prediction framework for the railway domain by combining XGBoost, a highly accurate tree-based model, with SHAP, a widely used explainability technique.

5- A Margin Based Early-Stopping Approach for Random Forest Classifiers

Status: Accepted and Published

DOI: <https://doi.org/10.1109/DSN-W65791.2025.00050>

Venue: IEEE/IFIP International Conference on Dependable Systems and Networks, Approximate Computing Workshop

Authors: Mario Barbareschi, Antonio Emmanuele

Indexes: Scopus

Abstract: In recent years, the continuous expansion of the Internet of Things (IoT) has driven a shift toward decentralized computation to improve privacy and reduce latency. As a result, computational tasks are increasingly executed directly on edge nodes within the IoT infrastructure. This includes the deployment of various Machine Learning models, such as Random Forest (RF) classifiers, which are widely used for their efficiency and interpretability. However, edge nodes are typically constrained in terms of computational power and memory, making it difficult to achieve high throughput when running complex RF models. Additionally, the large size of these models leads to higher energy consumption during inference, which can significantly impact the battery life of edge devices. To address these challenges, we propose an early-stopping mechanism for RF classifiers. The method consists of an offline phase, where statistical measures are precomputed to estimate the confidence level of predictions, and an online phase, where inference can be stopped early based on the computed confidence.

6- Exploiting Modular Redundancy for approximating Random Forest classifiers

Status: Under Review

Venue: Future Generation Computer Systems, Elsevier

Authors: Antonio Emmanuele, Mario Barbareschi, Alberto Bosio

Indexes: Scopus, WoS

Abstract: The deployment of machine learning models at the edge is crucial for enabling low-latency decision-making, optimizing resource utilization, and enhancing data confidentiality. Random Forest classifiers have proven to be highly accurate while offering computationally efficient inference, making them well-suited for resource-constrained edge devices. However, as the volume of training data grows, the complexity and size of these models also increase, limiting their deployment in edge computing scenarios.

In order to address this challenge, we propose a novel approximation strategy for Random Forest classifiers leveraging on the concept of modular redundancy.

In particular, our approach imposes that each target class is determined by only a subset of trees in a modular redundant fashion. This allows to prune from each tree the leaves related to no-longer relevant classes, significantly reducing the size of the model.

To achieve an optimal balance between accuracy and resource savings with minimal computational time, we introduce an heuristic algorithm that determine the best subset of trees for each class. We evaluate our approach on multiple UCI machine learning datasets using a hardware accelerator for tree ensembles, demonstrating its effectiveness. The result shows that, on average, a 2.5 % reduction in accuracy leads to save up to 50 % in hardware overhead and energy consumption.

7- Reliability analysis of hardware accelerators for decision tree-based classifier systems

Status: Under Review

Venue: Future Generation Computer Systems, Elsevier

Authors: Mario Barbareschi, Salvatore Barone, Antonio Emmanuele, Alberto Bosio

Indexes: Scopus, WoS

Abstract: The increasing adoption of AI models has driven applications toward the use of hardware accelerators to meet high computational demands and strict performance requirements. Beyond consideration of performance and energy efficiency, explainability and reliability have emerged as pivotal requirements, particularly for critical applications such as automotive, medical, and aerospace systems. Among the various AI models, Decision Tree Ensembles (DTEs) are particularly notable for their high accuracy and explainability. Moreover, they are particularly well-suited for hardware implementations, enabling high-performance and improved energy efficiency. However, a frequently overlooked aspect of DTEs is their reliability in the presence of hardware malfunctions. While DTEs are generally regarded as robust by design, due to their redundancy and voting mechanisms, hardware faults can still have catastrophic consequences. To address this gap, we present an in-depth reliability analysis of two types of DTE hardware accelerators: classical and approximate

implementations. Specifically, we conduct a comprehensive fault injection campaign, varying the number of trees involved in the classification task, the approximation technique used, and the tolerated accuracy loss, while evaluating several benchmark datasets. The results of this study demonstrate that designers must exercise caution when applying approximation techniques, as they can significantly impact resilience. Conversely, techniques that target the representation of features and thresholds appear to be better suited for fault tolerance.

8- vPUF: Virtualizing the Physical Unclonable Function

Status: Rejected after First Review Round

Venue: Original: IEEE Transaction on Emerging Topics in Computing

Authors: Mario Barbareschi, Antonio Emmanuele, Daniele Lombardi

Indexes: Scopus, WoS

Abstract: The adoption of Physical Unclonable Functions (PUFs) has become a prominent lightweight solution for enabling device authentication and secure key generation in resource-constrained IoT nodes. However, modern edge computing scenarios, where multiple isolated applications may coexist on the same physical device, limit the actual applicability of PUF-based security mechanisms since these generally assume the PUF as a primitive accessed by a single user. As a matter of fact, directly sharing a PUF among different users inherently undermines security-guarantees of these mechanisms. In addition, it jeopardizes isolation, as applications can impersonate each other by simply challenging the same PUF.

In this paper, we address these limitations by introducing the concept of the virtual PUF (vPUF): a novel abstraction that enables the secure and isolated sharing of a single PUF across multiple applications. We formally define the vPUF model, extend the set of classical PUF properties, and introduce new metrics to evaluate the quality and security of vPUF instances. Additionally, we devise four different virtualization strategies, detailing fundamental security requirements of the hosting platform. Finally, we prototype our solution on a RISC-V core running the Xvisor hypervisor and prove its effectiveness by directly measuring the discussed quality metrics, as well as additional required overhead.

9- A Dual-Stage Approach for Explainable Time Series Forecasting in On-Field Cyber-Physical Systems

Status: Under Review

Venue: Internet Of Things, Elsevier

Authors: Mario Barbareschi, Antonio Emmanuele, Nicola Mazzocca, Franca Rocco Di Torrepadula

Training and Research Activities Report

PhD in Information Technology and Electrical Engineering

Cycle: XXXIX

Author: Antonio Emmanuele

Indexes: Scopus, WoS

Abstract: Time series forecasting (TSF) plays a key role in Cyber-Physical Systems (CPS), enabling critical tasks, such as predictive maintenance, resource optimization, and anomaly detection. Recently, Deep Learning (DL) models have been widely employed for TSF, given their excellent capability in capturing temporal data relationships.

However, the computational demands of these models typically necessitate cloud-based deployment, which raises significant privacy and energy concerns. Additionally, the inherent complexity of these “black-box” models poses interpretability challenges, making it difficult for practitioners to understand and trust their decision-making processes.

To overcome these limitations, this paper proposes a dual-stage time series forecasting approach tailored for CPS applications. The approach is designed to meet the dual requirements of local execution at the edge and model interpretability. Specifically, in the first stage, clustering is applied to group time series data with similar patterns, while in the second stage, a Decision Tree Regressor (DTR) is trained for each cluster. By focusing on a narrower subset of data, each tree operates with reduced complexity, enabling the creation of simpler and more efficient models. Additionally, DTRs are inherently interpretable, offering a clear rationale for each prediction.

The lightweight, interpretable design supports efficient on-device execution with low computational demands, making it ideal for CPS scenarios. Experimental results on 250 time series demonstrate a significant improvement in accuracy compared to classical DTRs, while substantially improving interpretability and reducing computational time against tree-based ensembles, as validated on a Raspberry Pi 4 device.

5. Conferences and seminars attended

- 1- International Conference on Advanced Information Networking and Applications (AINA), 9-11 April Barcelona Spain (remote attendance), Paper Presented
- 2- European Dependable Computing Conference(EDCC), 8-11 April 2025, Lisbon Portugal, Paper Presented
- 3- IEEE/IFIP International Conference on Dependable Systems and Networks (DSN), 23-26 June, Naples Italy, Paper Presented (AxC workshop)
- 4- International Symposium on Design and Diagnostics of Electronic Circuits and Systems (DDECS), May 7-9, Lyon, France
- 5- International Conference on Availability, Reliability and Security (ARES), 11-14 August, Belgium, Ghent, 2025
- 6- 20th TAROT Summer School on Software Testing, Verification & Validation, June 30 – July 4 2025, Naples, Italy
- 7- IWES Ph.D. School in Embedded Systems, September 15 – September 17 2025, Modena, Italy

6. Periods abroad and/or in international research institutions

Period: 07/10/2025 – 31/10/2025

Supervisor: prof. Alberto Bosio

Hosting Institution : Institut des Nanotechnologies de Lyon, École Centrale de Lyon

Training and Research Activities Report

PhD in Information Technology and Electrical Engineering

Cycle: XXXIX

Author: Antonio Emmanuele

Research Activities: My research activities are aligned with the edge execution of machine learning models. Specifically, I am working on the following topics:

1. In-memory accelerator for Decision Tree models using Ternary Content Addressable Memories.
2. Investigation of the representativeness of faults in machine learning models.
3. Spatial architectures for accelerating Decision Trees.

In depth details regarding these works are provided in Section 8.

7. Tutorship

- 1- Architettura e Progetto dei Calcolatori (prof. Nicola Mazzocca), Corso di Laurea Magistrale in Ingegneria Informatica, hours: 30
- 2- Calcolatori Elettronici (prof. Mario Barbareschi), Corsi di Laurea triennale in Ingegneria Elettronica e Biomedica, hours: 10

8. Plan for year three

Research Activities and Period Abroad.

During my research activities, mostly carried out during my abroad period starting from the third year, I plan to explore the following topics related to the execution of ML models:

1. In-Memory Acceleration for Decision Trees.

In-Memory Computing is a technique aimed at reducing the computational bottleneck caused by memory transfers, which increase latency due to memory stalls and energy consumption due to frequent accesses to higher-level memories. In this context, In-Memory Computing enables performing computation directly near memory elements, thus avoiding unnecessary data transfers. The main objective is to accelerate Decision Tree inference by leveraging **Ternary Content Addressable Memories (TCAMs)**, memory cells with ternary values (0, 1, Don't Care), which are particularly useful for efficient rule matching during Decision Tree inference.

2. Spatial Architectures for Decision Tree Inference.

CMOS- and FPGA-based accelerators for Decision Trees can be categorized according to their nature. High-throughput accelerators are typically model-specific, offering minimal inference latency but limited flexibility. Conversely, non-model-specific accelerators usually perform generic tree traversal, providing higher flexibility and throughput at the cost of increased latency. Spatial architectures—such as the well-known **Systolic Array**—represent a promising trade-off between flexibility and latency. On one hand, similar to non-model-specific accelerators, the model can still be represented as a memory element. On the other hand, by performing computation in space rather than in time, it becomes possible to explore multiple trees and multiple levels within the same tree simultaneously.

3. Fault Representativeness.

While Fault Injection techniques are effective for analyzing how hardware faults impact machine learning models in terms of accuracy, they still lack the ability to provide meaningful insight into the contribution of each fault to a specific prediction. Understanding the individual contribution of a fault during inference would enable identifying which model components are most sensitive to faults.

Training and Research Activities Report

PhD in Information Technology and Electrical Engineering

Cycle: XXXIX

Author: Antonio Emmanuele

This understanding could guide ML model developers to design models that account for component susceptibility to faults, while hardware designers could reinforce accelerator components responsible for the computation of more critical model parts during prediction.

Thesis Proposal

In light of the activities already completed and those being carried out during my period abroad, my thesis will encompass the key aspects essential for the adoption of Decision Tree models at the edge:

1. **Hardware Acceleration on Edge Architectures.**

This part will address all aspects of Decision Tree acceleration explored during the three years, including implementations on microcontrollers, FPGAs (spatial architectures), and In-Memory Computing architectures.

2. **Approximation of Decision Tree-Based Models.**

This section will follow the hardware acceleration part and discuss approximation techniques for machine learning models, which are essential for deploying compressed models with reduced computational, memory, and hardware requirements.

3. **Fault Resiliency of Decision Tree Models.**

Finally, to enable the full adoption of Decision Tree models in safety-critical contexts, it is essential to investigate their fault reliability, providing insights into the behaviour of different accelerators and approximation techniques under fault conditions.