Joint ICTP-IAEA School on Detector Signal Processing and Machine Learning for Scientific Instrumentation and Reconfigurable Computing

Introduction to Machine Learning and Edge Al

smr4110 | Trieste, Italy - 2025

Romina Soledad Molina, Ph.D.



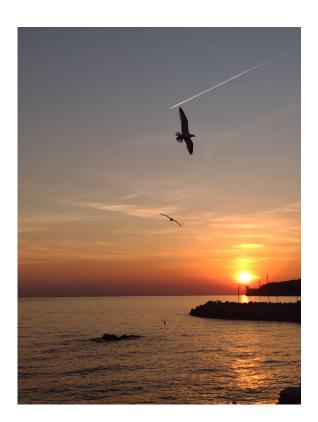






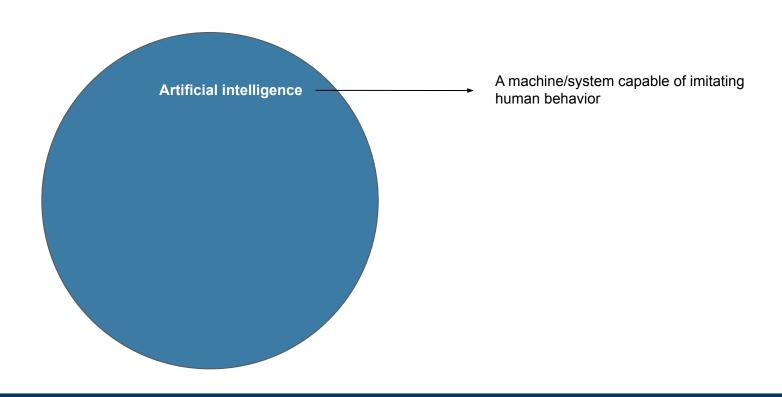
Outline

- Introduction.
- Al beyond 2025.
- What is Edge AI?
- What the Literature Says.
- Integrating Every Layer of Edge AI Design.
- Case Study for FPGA-based Edge AI: MNIST Binary Classification.

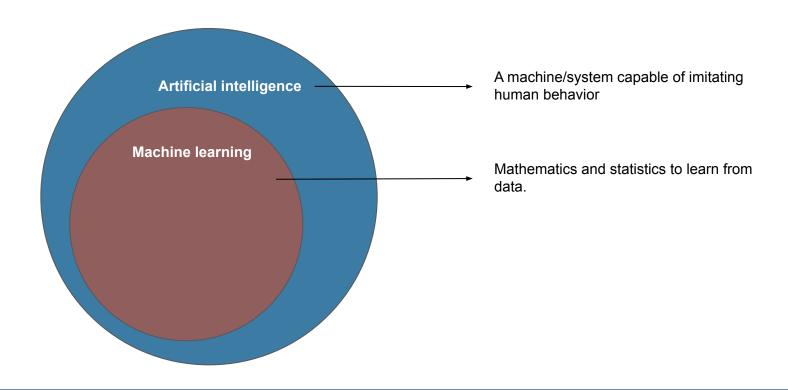




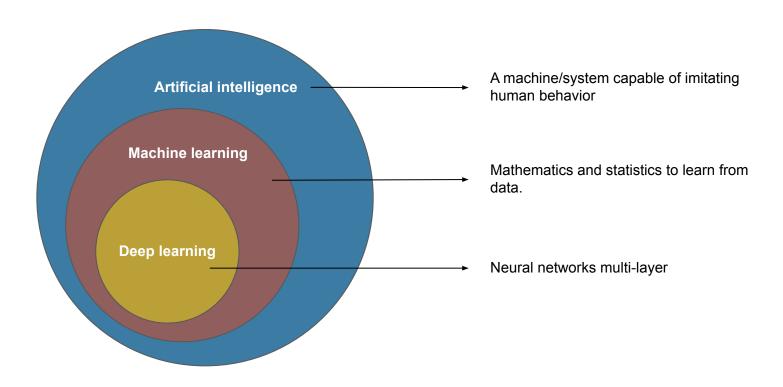




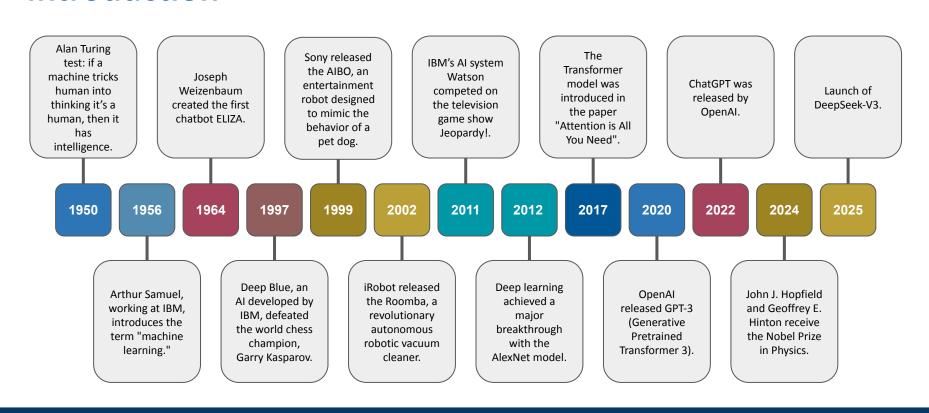














Why has the explosion of AI happened in recent years? Three main components

Big data



Why has the explosion of AI happened in recent years? Three main components

Big data

Software



Why has the explosion of AI happened in recent years? Three main components

Big data

Software

Hardware



Why has the explosion of AI happened in recent years? Three main components

Big data

Software

Hardware

The massive amount of data generated every day—from social media, sensors, and online activity—provides the fuel that AI systems need to learn and improve



Why has the explosion of AI happened in recent years? Three main components

Big data

The massive amount of data generated every day—from social media, sensors, and online activity—provides the fuel that AI systems need to learn and improve

Software

New algorithms, neural network architectures, and open-source frameworks have dramatically improved Al's capabilities and made development more accessible. Hardware



Why has the explosion of AI happened in recent years? Three main components

Big data

The massive amount of data generated every day—from social media, sensors, and online activity—provides the fuel that AI systems need to learn and improve

Software

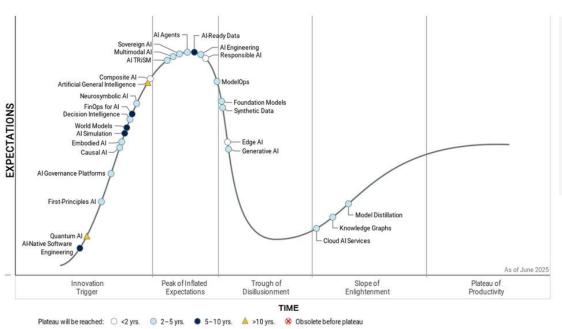
New algorithms, neural network architectures, and open-source frameworks have dramatically improved Al's capabilities and made development more accessible.

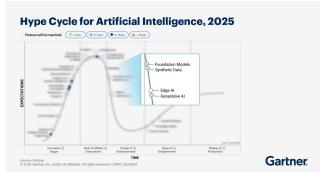
Hardware

Advances in GPUs, TPUs, and cloud computing have made it possible to train large AI models faster and more efficiently than ever before.



AI beyond 2025





Graphically shows the maturity and adoption of technologies, helping to solve real problems and take advantage of opportunities.

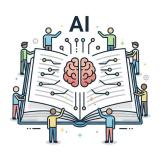
Gartner.

Source: https://www.gartner.com/en/newsroom/press-releases/2025-08-05-gartner-hype-cycle-identifies-top-ai-innovations-in-2025





Al Democratization: Making Al tools and knowledge accessible to everyone, lowering barriers for innovation and participation worldwide.







AI beyond 2025

Al Democratization: Making Al tools and knowledge accessible to everyone, lowering barriers for innovation and participation worldwide.

Ethical AI and Regulation: Promoting transparency, fairness, and accountability through responsible AI development and global policy frameworks.







Al Democratization: Making Al tools and knowledge accessible to everyone, lowering barriers for innovation and participation worldwide.

Ethical AI and Regulation: Promoting transparency, fairness, and accountability through responsible AI development and global policy frameworks.

AI in Healthcare: Powering precision medicine, early diagnostics, and improved patient care through data-driven insights.







Al Democratization: Making Al tools and knowledge accessible to everyone, lowering barriers for innovation and participation worldwide.

Ethical AI and Regulation: Promoting transparency, fairness, and accountability through responsible AI development and global policy frameworks.

Al in Healthcare: Powering precision medicine, early diagnostics, and improved patient care through data-driven insights.

Al and Human Collaboration: Enhancing creativity and productivity by combining human intuition with machine intelligence.







Al Democratization: Making Al tools and knowledge accessible to everyone, lowering barriers for innovation and participation worldwide.

Ethical AI and Regulation: Promoting transparency, fairness, and accountability through responsible AI development and global policy frameworks.

Al in Healthcare: Powering precision medicine, early diagnostics, and improved patient care through data-driven insights.

Al and Human Collaboration: Enhancing creativity and productivity by combining human intuition with machine intelligence.

Autonomous Systems and Al-Powered Robotics: Enabling intelligent automation in transportation, manufacturing, and logistics.







Al in Education: Delivering personalized and adaptive learning experiences for students across all levels.



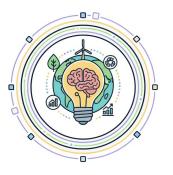




Al in Education: Delivering personalized and adaptive learning experiences for students across all levels.

Natural Language Processing (NLP) and Conversational AI: Advancing communication between humans and machines through natural, context-aware interactions.







Al in Education: Delivering personalized and adaptive learning experiences for students across all levels.

Natural Language Processing (NLP) and Conversational AI: Advancing communication between humans and machines through natural, context-aware interactions.

Quantum Al: Leveraging quantum computing to accelerate learning and solve complex problems beyond classical limits.







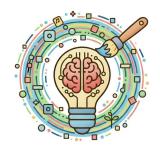
AI beyond 2025

Al in Education: Delivering personalized and adaptive learning experiences for students across all levels.

Natural Language Processing (NLP) and Conversational AI: Advancing communication between humans and machines through natural, context-aware interactions.

Quantum AI: Leveraging quantum computing to accelerate learning and solve complex problems beyond classical limits.

Al and Sustainability: Applying Al to address climate change, optimize energy use, and build greener technologies.







Al in Education: Delivering personalized and adaptive learning experiences for students across all levels.

Natural Language Processing (NLP) and Conversational AI: Advancing communication between humans and machines through natural, context-aware interactions.

Quantum AI: Leveraging quantum computing to accelerate learning and solve complex problems beyond classical limits.

Al and Sustainability: Applying Al to address climate change, optimize energy use, and build greener technologies.

Al-Driven Creativity: Unlocking new forms of art, design, and innovation through generative models.







5G/6G and Edge AI: Combining ultra-fast connectivity with local AI processing to enable real-time, low-latency intelligence at the edge.







AI beyond 2025

5G/6G and Edge AI: Combining ultra-fast connectivity with local AI processing to enable real-time, low-latency intelligence at the edge.

Al Infrastructure and MLOps: Building scalable pipelines for data, model training, and deployment to support enterprise Al adoption.







5G/6G and Edge AI: Combining ultra-fast connectivity with local AI processing to enable real-time, low-latency intelligence at the edge.

Al Infrastructure and MLOps: Building scalable pipelines for data, model training, and deployment to support enterprise Al adoption.

Al and Cybersecurity: Strengthening digital defense systems and detecting threats proactively using intelligent algorithms.







AI beyond 2025

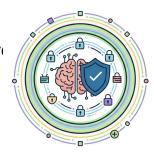
5G/6G and Edge AI: Combining ultra-fast connectivity with local AI processing to enable real-time, low-latency intelligence at the edge.

Al Infrastructure and MLOps: Building scalable pipelines for data, model training, and deployment to support enterprise Al adoption.

Al and Cybersecurity: Strengthening digital defense systems and detecting threats proactively using intelligent algorithms.

Neuromorphic Computing: Designing brain-inspired chips for efficient, adaptive AI.







Question: What do you think about the future of AI?





Al Market Growth and Impact

Market Growth:

 The global AI market is currently valued at several hundred billion USD and is projected to reach between US \$2–4 trillion by the early-2030s, with annual growth rates (CAGR) in the ≈ 20% to 31% range — depending on the forecast.



Al Market Growth and Impact

Market Growth:

 The global AI market is currently valued at several hundred billion USD and is projected to reach between US \$2–4 trillion by the early-2030s, with annual growth rates (CAGR) in the ≈ 20% to 31% range — depending on the forecast.

Business Applications:

- According to Gartner, Inc., by 2026 up to 40% of enterprise applications will contain task-specific Al agents (up from less than 5% today).
- They further predict that by 2029, 80% of common customer-service issues will be handled autonomously by what they call "agentic AI" (with little or no human intervention).



Al Market Growth and Impact

Market Growth:

 The global AI market is currently valued at several hundred billion USD and is projected to reach between US \$2–4 trillion by the early-2030s, with annual growth rates (CAGR) in the ≈ 20% to 31% range — depending on the forecast.

Business Applications:

- According to Gartner, Inc., by 2026 up to 40% of enterprise applications will contain task-specific Al agents (up from less than 5% today).
- They further predict that by 2029, 80% of common customer-service issues will be handled autonomously by what they call "agentic AI" (with little or no human intervention).

Impact on Businesses & Jobs:

 Business models and job roles are poised to undergo rapid and sometimes unpredictable transformations—driven by autonomous systems, embedded AI in workflows and changing expectations of productivity and capabilities.

AI-Specific Hardware Platforms

- Al Chips & Processors
 - GPUs (e.g., NVIDIA H100, AMD Instinct MI300) Still the backbone of Al acceleration.
 - TPUs (Google Tensor Processing Units) Optimized for deep learning workloads.
 - NPUs (Neural Processing Units) Specialized for on-device AI in smartphones and edge devices.
 - Analog & Optical Al Chips Emerging technology for ultra-low-power Al inference.



AI-Specific Hardware Platforms

- Al Hardware Trends
 - Edge AI Custom low-power AI chips for real-time processing on IoT devices.
 - RISC-V AI Processors Open-source architecture gaining traction for AI acceleration.
 - Neuromorphic Computing Brain-inspired hardware and algorithms designed for parallel, adaptive, and energy-efficient processing.



AI-Specific Hardware Platforms

Question: What do you think about FPGA and AI?





AI-Specific Hardware Platforms

FPGA+Al Trends

- AI-Specific FPGA Architectures New FPGA models optimized for deep learning (Xilinx Versal AI Core, Intel Stratix 10 NX, and Lattice Avant AI).
- Quantum + FPGA Hybrid Computing Research into integrating FPGA with quantum acceleration.
- FPGA as-a-Service (FaaS) Cloud-based FPGA solutions for scalable Al workloads.
- Embedded AI & Edge Computing FPGAs in IoT & industrial AI, enabling real-time decision-making with ultra-low power.





"Edge artificial intelligence (Edge AI) refers to the deployment of AI algorithms and AI models directly on local edge devices such as sensors or Internet of Things (IoT) devices, which enables real-time data processing and analysis without constant reliance on cloud infrastructure." [IBM]



"Edge artificial intelligence (Edge AI) refers to the deployment of AI algorithms and AI models directly on local edge devices such as sensors or Internet of Things (IoT) devices, which enables real-time data processing and analysis without constant reliance on cloud infrastructure." [IBM]

On the edge processing



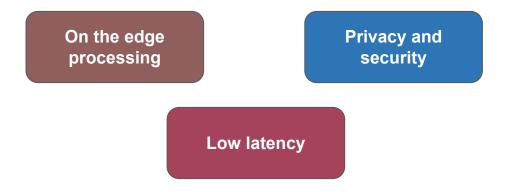
"Edge artificial intelligence (Edge AI) refers to the deployment of AI algorithms and AI models directly on local edge devices such as sensors or Internet of Things (IoT) devices, which enables real-time data processing and analysis without constant reliance on cloud infrastructure." [IBM]

On the edge processing

Low latency



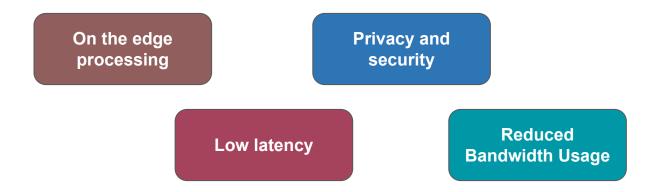
"Edge artificial intelligence (Edge AI) refers to the deployment of AI algorithms and AI models directly on local edge devices such as sensors or Internet of Things (IoT) devices, which enables real-time data processing and analysis without constant reliance on cloud infrastructure." [IBM]



[IBM] https://www.ibm.com/think/topics/edge-ai



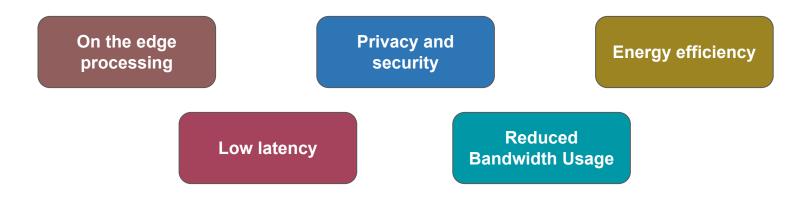
"Edge artificial intelligence (Edge AI) refers to the deployment of AI algorithms and AI models directly on local edge devices such as sensors or Internet of Things (IoT) devices, which enables real-time data processing and analysis without constant reliance on cloud infrastructure." [IBM]



[IBM] https://www.ibm.com/think/topics/edge-ai



"Edge artificial intelligence (Edge AI) refers to the deployment of AI algorithms and AI models directly on local edge devices such as sensors or Internet of Things (IoT) devices, which enables real-time data processing and analysis without constant reliance on cloud infrastructure." [IBM]



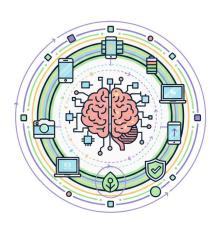
[IBM] https://www.ibm.com/think/topics/edge-ai



Intelligence closer to data

- Low latency: Real-time decisions, reducing cloud dependence.
- Privacy & security: Data stays on the device.
- Lower bandwidth use: Less data sent to the cloud.
- Energy efficiency: Optimized local processing.
- Scalability & reliability: Devices run autonomously, even offline

Bringing intelligence closer to the data, enabling low-latency, privacy-preserving, and energy-aware Al.





Edge Al Challenges

Hardware & Energy

Computation & Performance

Data & Deployment

Limited power

Latency & real-time

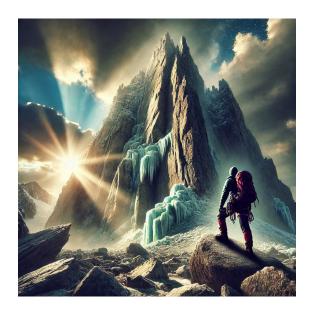
Privacy & security

Memory/storage limits

Limited compute capacity

Network issues

Hardware heterogeneity





Edge Al Challenges

• Question:

Given all these challenges... how can we make AI models run efficiently on tiny, power-limited devices?





Edge Al Challenges

• Question:

What if we could combine the strengths of different processors CPUs, GPUs, and FPGAs — to overcome these limits?





Heterogeneous Computing: A Key Enabler for Edge AI.

By incorporating CPUs, GPUs, and FPGAs into a single system,
 heterogeneous computing becomes possible.



Heterogeneous Computing: A Key Enabler for Edge AI.

- By incorporating CPUs, GPUs, and FPGAs into a single system,
 heterogeneous computing becomes possible.
- This approach maximizes the strengths of each component, efficiently distributing edge workloads to improve both performance and energy consumption.



Heterogeneous Computing: A Key Enabler for Edge AI.

• **CPU:** handles general-purpose tasks like system management, control logic, and lightweight AI inference.



Heterogeneous Computing: A Key Enabler for Edge AI.

- **CPU:** handles general-purpose tasks like system management, control logic, and lightweight AI inference.
- **GPU:** accelerates parallelizable tasks such as deep learning inference, image processing, and high-performance computing.



Heterogeneous Computing: A Key Enabler for Edge AI.

- **CPU:** handles general-purpose tasks like system management, control logic, and lightweight AI inference.
- GPU: accelerates parallelizable tasks such as deep learning inference, image processing, and high-performance computing.
- **FPGA:** provides real-time, ultra-low-latency processing for tasks like sensor fusion, encryption, and custom AI accelerators.





Low latency



Low latency

Energy Efficiency



Low latency

Energy Efficiency

High parallelism



Low latency

Energy Efficiency

High parallelism

Scalability



Low latency

Energy Efficiency

High parallelism

Scalability

Customizable Al Acceleration



Low latency

Energy Efficiency

High parallelism

Scalability

Customizable Al Acceleration

FPGA / SoC-based on FPGA

Resource-constrained devices



Original ML-based model

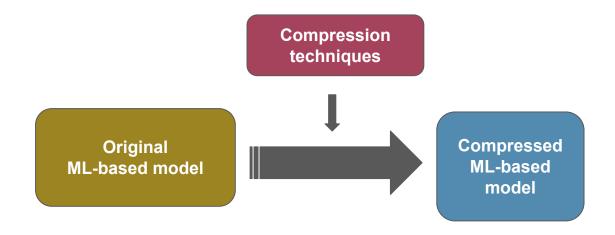


Original ML-based model

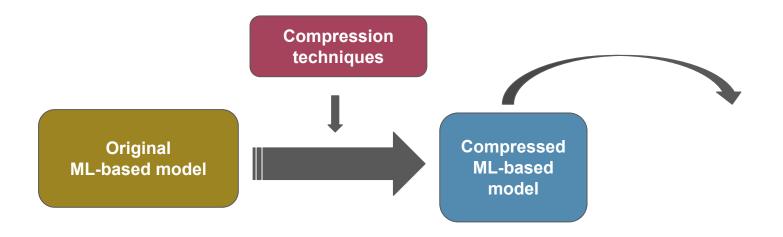




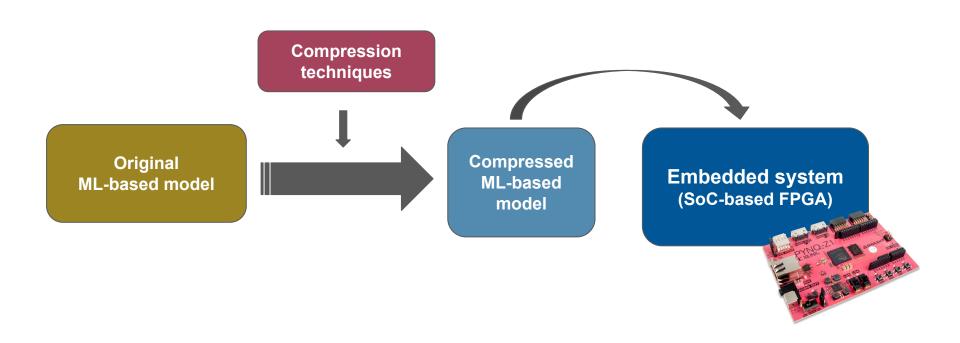
















- ML on constrained devices is now a major focus
 - TinyML and on-device inference are growing research areas.
- Classical state-of-the-art models are being adapted for the edge
 - MobileNetV2, TinyBERT, and YOLO-Lite show that efficiency is now as important as accuracy.
- Compression through pruning & quantization dominates optimization
 - These remain the most used techniques to fit large models into edge hardware.
- Off-chip memory access is a key bottleneck
 - Memory transfers can consume more energy than computation.
- Workflows are still fragmented across tools
 - Full end-to-end edge pipelines are rare; integration remains an open challenge.



Memory footprint and latency



Memory footprint and latency

Ensemble of compression techniques



Memory footprint and latency

Ensemble of compression techniques

On-chip memory deployment



What the Literature Says

Memory footprint and latency

Ensemble of compression techniques

On-chip memory deployment

End-to-end workflow



What the Literature Says

Memory footprint and latency

Ensemble of compression techniques

On-chip memory deployment

End-to-end workflow

Productivity





Optimization is holistic — every phase, from model to hardware, impacts overall performance at the edge.

Software development

Model training

Model optimization & compression



Optimization is holistic — every phase, from model to hardware, impacts overall performance at the edge.

Software development

Model training

Model optimization & compression

Hardware platform selection

Edge devices

Cloud/GPU/CPU

Other computing platforms



Optimization is holistic — every phase, from model to hardware, impacts overall performance at the edge.

Software development

Model training

Model optimization & compression

Hardware platform selection

Edge devices

Cloud/GPU/CPU

Other computing platforms

Firmware creation

C/C++ application

Python application

HDL design



Optimization is holistic — every phase, from model to hardware, impacts overall performance at the edge.

Software development

Model training

Model optimization & compression

Hardware platform selection

Edge devices

Cloud/GPU/CPU

Other computing platforms

Firmware creation

C/C++ application

Python application

HDL design

PCB design

Key components

Signal and power integrity

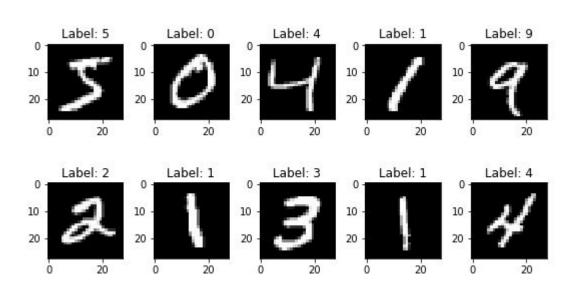
Methodological design

EMC and **EMI**

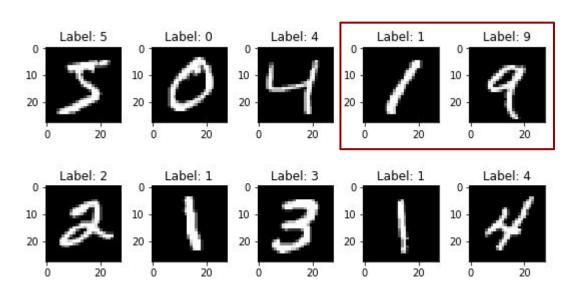


Case Study for FPGA-based Edge AI: MNIST Binary Classification.











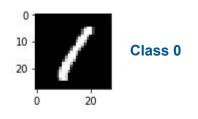
- Quantization-Aware Pruning
 - 8-bits fixed point precision
 - 20% target sparsity
 - QKeras for model definition

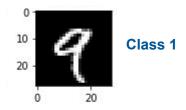
Layer (type)	Output	Shape	Param #
fc1_input (QDense)	(None,	5)	3925
relu_input (QActivation)	(None,	5)	0
fc1 (QDense)	(None,	7)	42
relu1 (QActivation)	(None,	7)	0
fc2 (QDense)	(None,	10)	80
relu2 (QActivation)	(None,	10)	Θ
output (QDense)	(None,	2)	22
sigmoid (Activation)	(None,	2)	0
======================================			

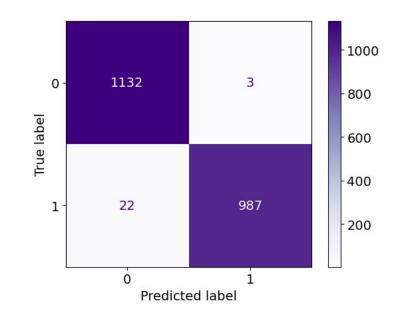


ML and model compression techniques for SoC/FPGA

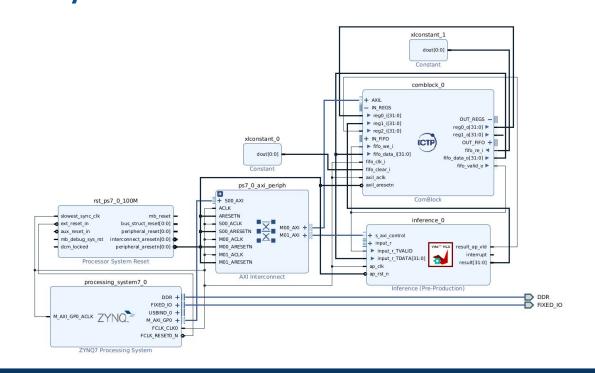
Demo: MNIST-based binary classification



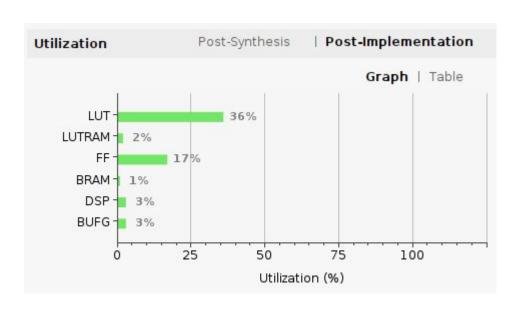




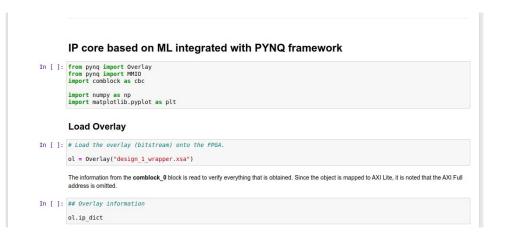
















ML and model compression techniques for SoC/FPGA

Demo: MNIST-based binary classification







```
Data preparation
In [ ]: signal 1 = [
        In []: imageArray = np.array(signal 1)
     image 2d = imageArray.reshape((28, 28))
     # Display as an image
     plt.imshow(image 2d, cmap='gray', interpolation='nearest')
     plt.colorbar() # Optional: Show color scale
     plt.show()
In [ ]: signal 2 = [
        In []: imageArray = np.array(signal 2)
     image 2d = imageArray.reshape((28, 28))
     # Display as an image
     plt.imshow(image 2d, cmap='gray', interpolation='nearest')
     plt.colorbar() # Optional: Show color scale
     plt.show()
```











Joint ICTP-IAEA School on Detector Signal Processing and Machine Learning for Scientific Instrumentation and Reconfigurable Computing

Introduction to Machine Learning and Edge Al

smr4110 | Trieste, Italy - 2025

Romina Soledad Molina, Ph.D.



