# Hardware Architectures for Embedded and Edge AI (from ML to HW and back)
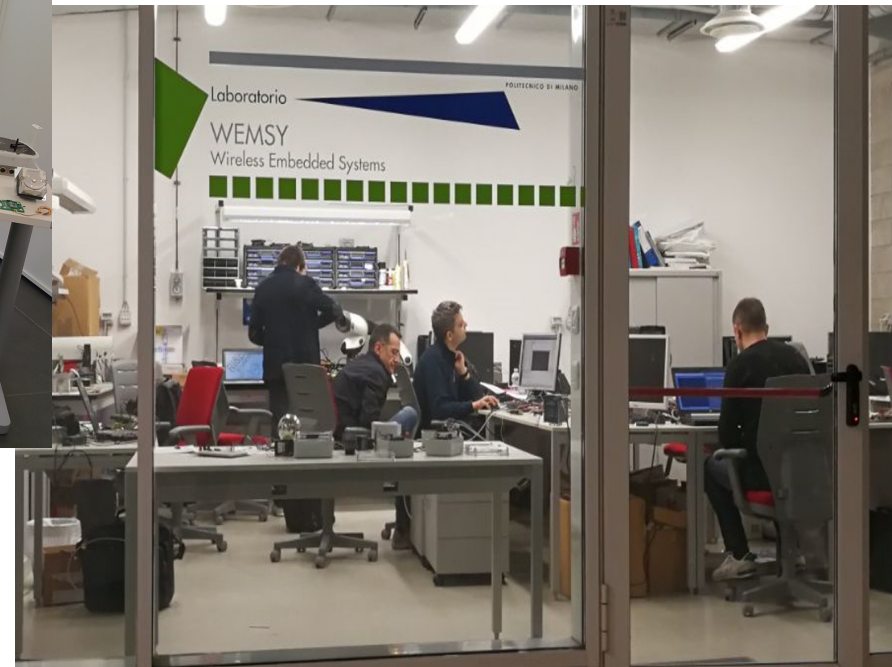
*Prof. Manuel Roveri*

*«Workshop on Widening Access to TinyML Network by Establishing Best Practices in Education»*

# Prof. Manuel Roveri



- **Full Professor**

  Dipartimento di Elettronica, Informazione e Bioingegneria (DEIB), Politecnico di Milano, Italy
  **Email: manuel.roveri@polimi.it**
  Web: http://roveri.faculty.polimi.it

- **Research interests**: TinyML, IoT and edge computing, privacy-preserving machine and deep learning

- **Lecturer of « Computing Infrastructures» and «Hardware Architecture for Embedded and edge AI»**

- **Associate Editor** of IEEE Trans. on Artificial Intelligence, Neural Networks, IEEE Trans. on Emerging Tecnologies in Computational Intelligence, IEEE Trans. on Neural Networks and Learning Systems

- Chair of the IEEE CIS **Technical Activities** strategic planning committee and IEEE CIS **Neural Network** Technical Committee

- **Co-Founder of DHIRIA**, a Spin-Off of Politecnico di Milano

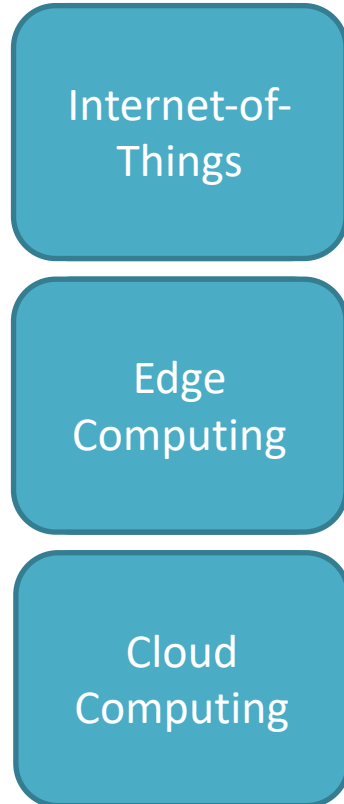# AI-Tech Research Lab @ Politecnico di Milano

# The research activity

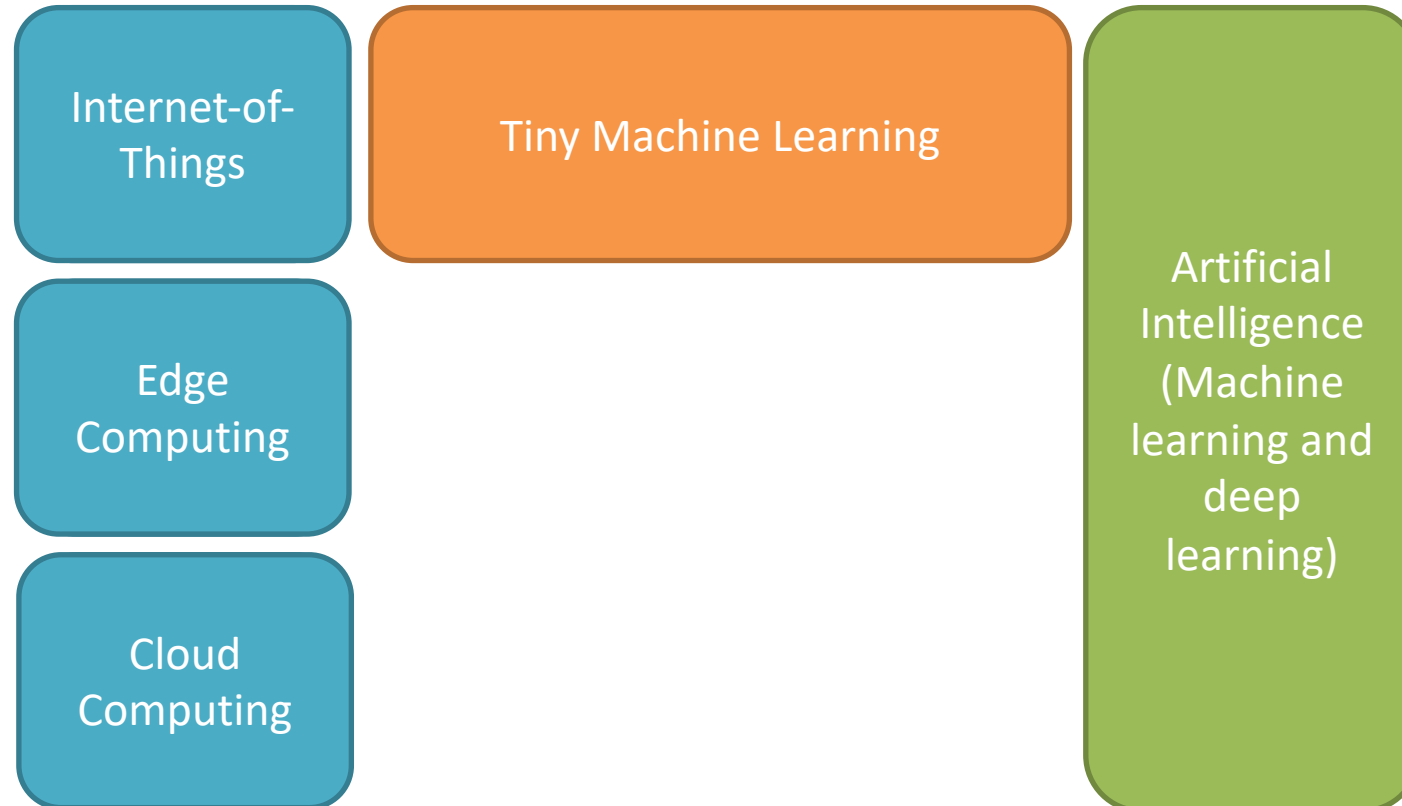Cyber-physical Systems

Artificial Intelligence (Machine learning and deep learning)

# The research activity

Internet-of-Things

Edge Computing

Cloud Computing

Artificial Intelligence (Machine learning and deep learning)

# The research activity

# The research activity

Internet-of-Things

Edge Computing

Cloud Computing

Distributed Inference and Federated Learning

Artificial Intelligence (Machine learning and deep learning)

# The research activity



Internet-of-Things

Edge Computing

Cloud Computing

Machine Learning as a service, Privacy-preserving Machine Learning

Artificial Intelligence (Machine learning and deep learning)

# The research activity



"Hardware Architectures for Embedded and Edge AI" Course

Internet-of-Things

Edge Computing

Cloud Computing

Tiny Machine Learning

Distributed Inference and Federated Learning

Machine-Learning-as-a-service, Privacy-preserving Machine Learning

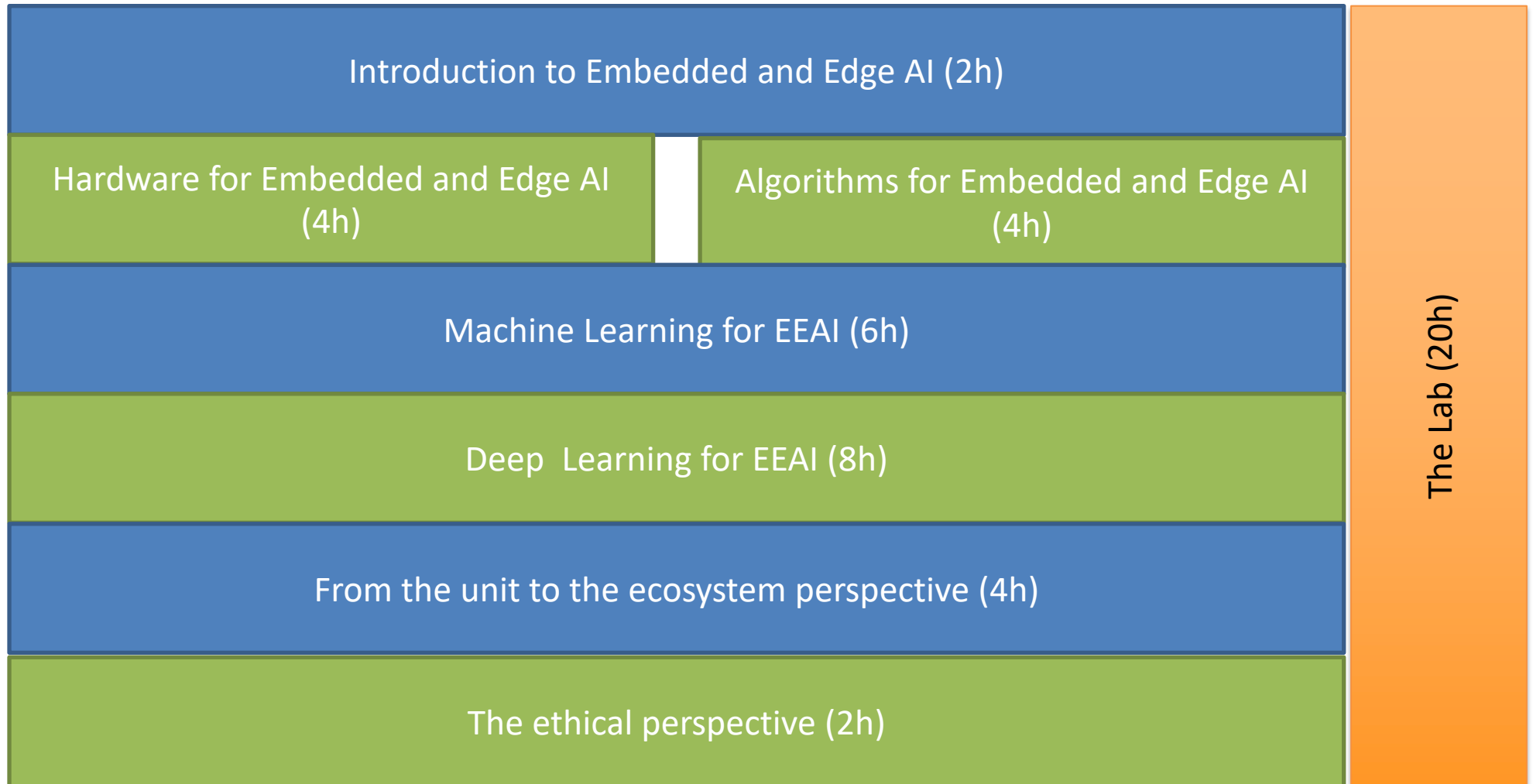Artificial Intelligence (Machine learning and deep learning)

"Computing Infrastructures" Course

"Hardware Architectures for Embedded and Edge AI"
Information about the course

# Course Details

- Course Title: "HARDWARE ARCHITECTURES FOR EMBEDDED AND EDGE AI"
- Academic Year  2022/2023
- School of Industrial and Information Engineering
- Master of Science degree - Computer Science and Engineering
- Course Type  Mono-Disciplinary Course
- Credits (CFU / ECTS)  5.0
- Course Organization: 30h lectures (M. Roveri) + 20h labs (M. Pavan)
- Number of enrolled students: 62
  - 66% Computer Science
  - 27% Electronics
  - 7% Bio – Control Theory - Telecom

# Course Organization



Introduction to Embedded and Edge AI (2h)

Hardware for Embedded and Edge AI (4h)

Algorithms for Embedded and Edge AI (4h)

Machine Learning for EEAI (6h)

Deep Learning for EEAI (8h)

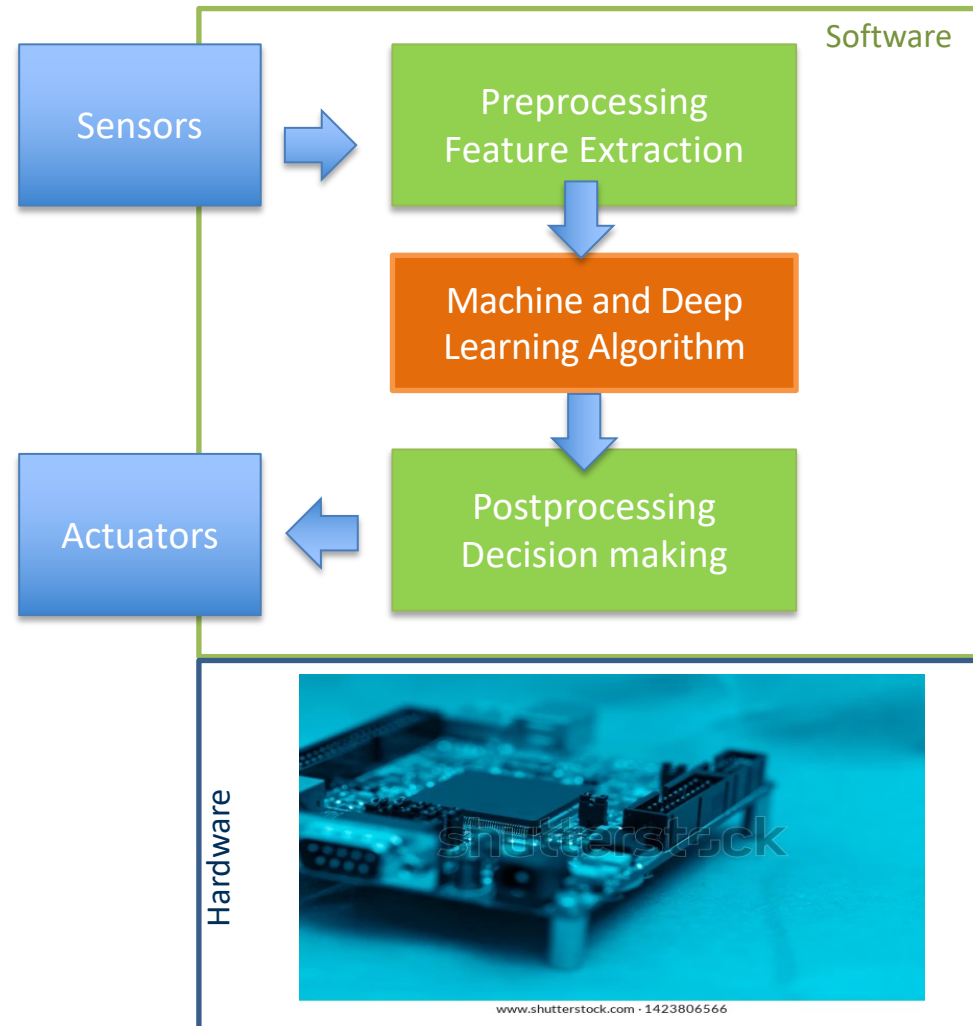From the unit to the ecosystem perspective (4h)

The ethical perspective (2h)

The Lab (20h)

# 1) Introduction to Embedded and Edge AI (2h)


Wake-word detection


Person detection


Gesture recognition

**Software**

Sensors → Preprocessing Feature Extraction → Machine and Deep Learning Algorithm → Postprocessing Decision making → Actuators

**Hardware**



Five Ws in Embedded and Edge Ai:
- Why do we need EEAI?
- What can we do with EEAI?
- Where can we find?
- When do we need it (design)?
- Who is in charge of EEAI code?
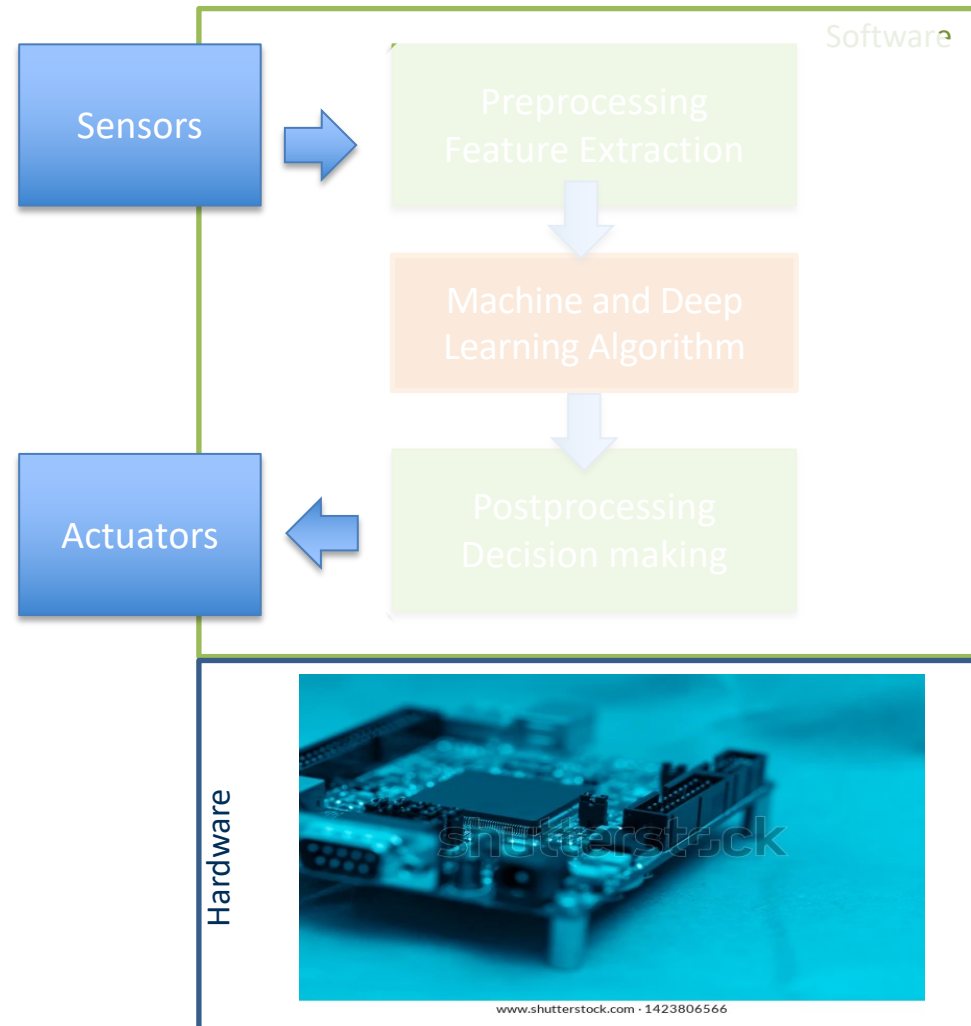
# 2) Hardware for Embedded and Edge AI (4h)

Wake-word detection

Person detection

Gesture recognition

Sensors → Preprocessing Feature Extraction

Software

Machine and Deep Learning Algorithm

Actuators ← Postprocessing Decision making

Hardware

- **Sensors and signals**: the TinyML perspective (ts, audio, image, video)
- **Sensors (+ application)**:
  - Acoustic and vibration
  - Visual and scene
  - Motion and position
  - Force and tactile
  - Optical and electromagnetic
  - Environmental and chemical

| MPUs | Low-end MCUs |
|------|-------------|
| High-end MCUs | SoCs |

Memory, computation, energy, cost

# 3) Algorithms for Embedded and Edge AI (4h)

Wake-word detection

Person detection

Gesture recognition

Software

Sensors

Preprocessing
Feature Extraction

Machine and Deep
Learning Algorithm

Postprocessing
Decision making

Actuators

Hardware

Chopping, Windowing

Reconstruction of missing data

Resampling

Filtering

Feature extraction

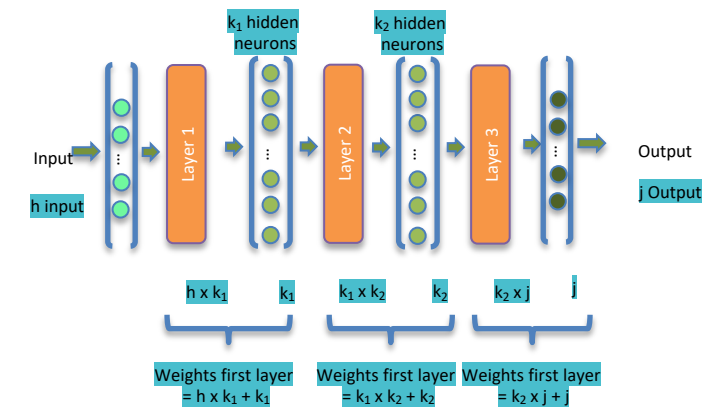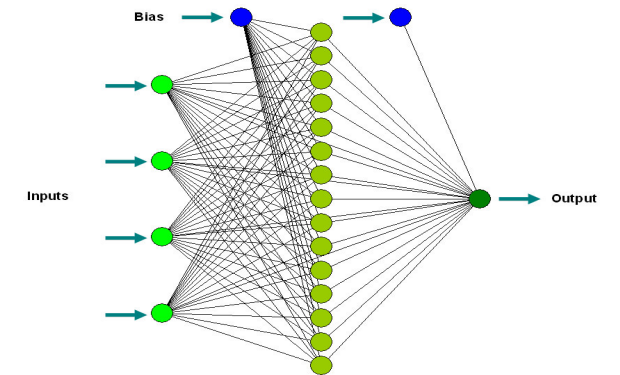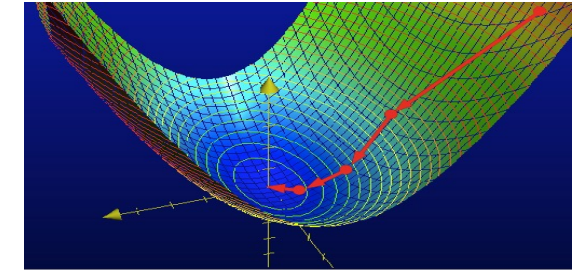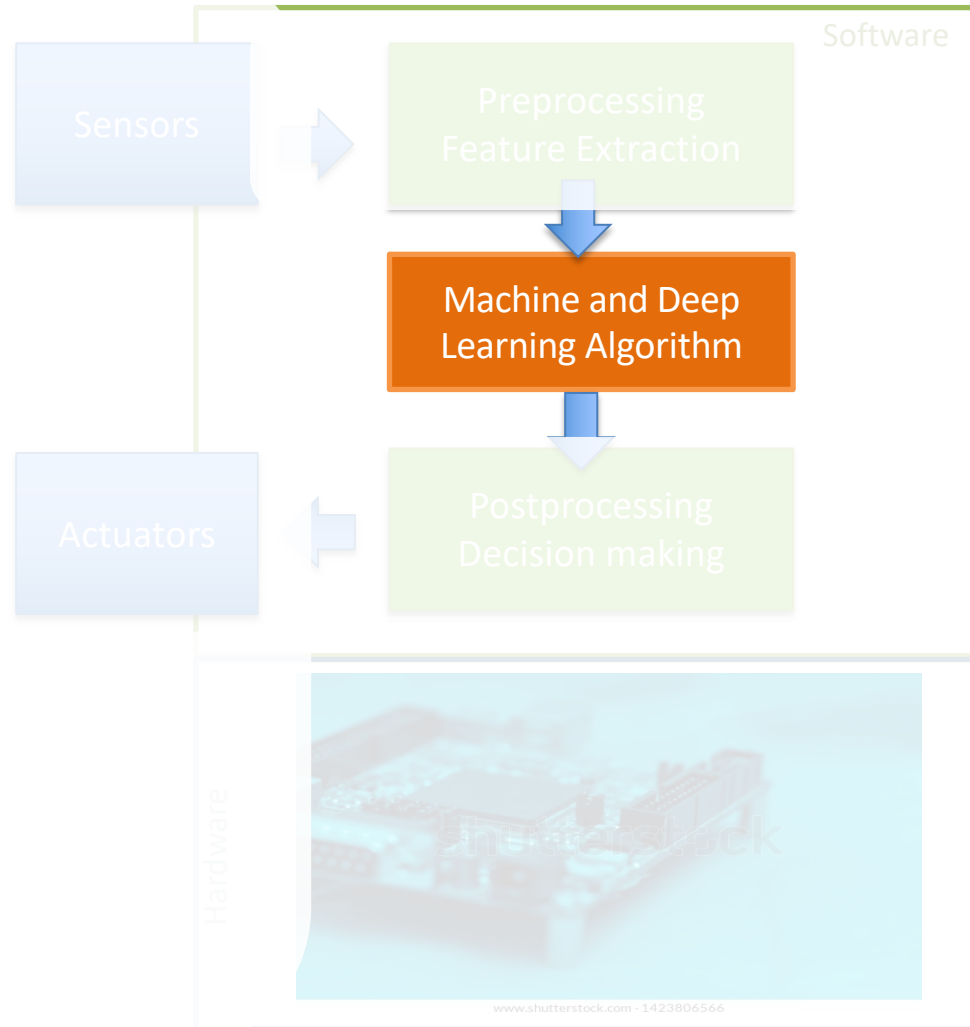# 4) Machine Learning for Embedded and Edge AI (6h)



Wake-word detection

Person detection

Gesture recognition

Software

Sensors

Preprocessing
Feature Extraction

Machine and Deep
Learning Algorithm

Actuators

Postprocessing
Decision making

Hardware

# 5) Deep Learning for Embedded and Edge AI (8h)



=≈ **37KB**

| 4KB | 19KB | 5KB | 6KB | 1.6KB | 0.5KB | 0.3KB | ...KB |

INPUT 32x32

C1: feature maps 6@28x28

S2: feature maps 6@14x14

C3: feature maps 16@10x10

S4: f. maps 16@5x5

C5: layer 120

F6: layer 84

OUTPUT 10

CONVOLUTIONS   SUBSAMPLING   CONVOLUTIONS   SUBSAMPLING   FULL CONNECTION   GAUSSIAN CONNECTIONS

6 5x5 Filters

16 5x5 Filters

FC1 400->120

FC2 120->84

**#Weights**

6x5x5x1+6 → **=156**

16x5x5x6+16 → **=2416**

400x120+120 → **=48120**

120x84+84 → **=10164**

**#MAC**

6x5x5x28x28 → **= 117K**

16x5x5x6x10x10 → **=240K**

400x120 → **=48K**

120x84 → **≈ 10K**

# 5) Deep Learning for Embedded and Edge AI (8h)

# 5) Deep Learning for Embedded and Edge AI (8h)



Architectures for EEAI → Approximate Computing → Embedded System Code Optimization

Output Layer

# 5) Deep Learning for Embedded and Edge AI (8h)

Architectures for EEAI ➡ Approximate Computing ➡ Embedded System Code Optimization

**SqueezeNet (2016)**

**MobileNet (2017)**

**EfficientNet (2019)**

Cit. Scholar 12/22: 7K

Cit. Scholar 12/22: 16K

Cit. Scholar 12/22: 9.5K

# 5) Deep Learning for Embedded and Edge AI (8h)

| Architectures for EEAI | ➡ | Approximate Computing | ➡ | Embedded System Code Optimization |
|---|---|---|---|---|

- Precision scaling:
  - ✓ **Quantization mechanisms**
  - ✓ **Implementation**
  - ✓ **Learning quantized models (PTQ, QAT)**
- Task dropping:
  - ✓ **network pruning**
  - ✓ **network architecture design**
  - ✓ **transfer learning**
  - ✓ **knowledge distillation**
- Early-exit Neural Networks:
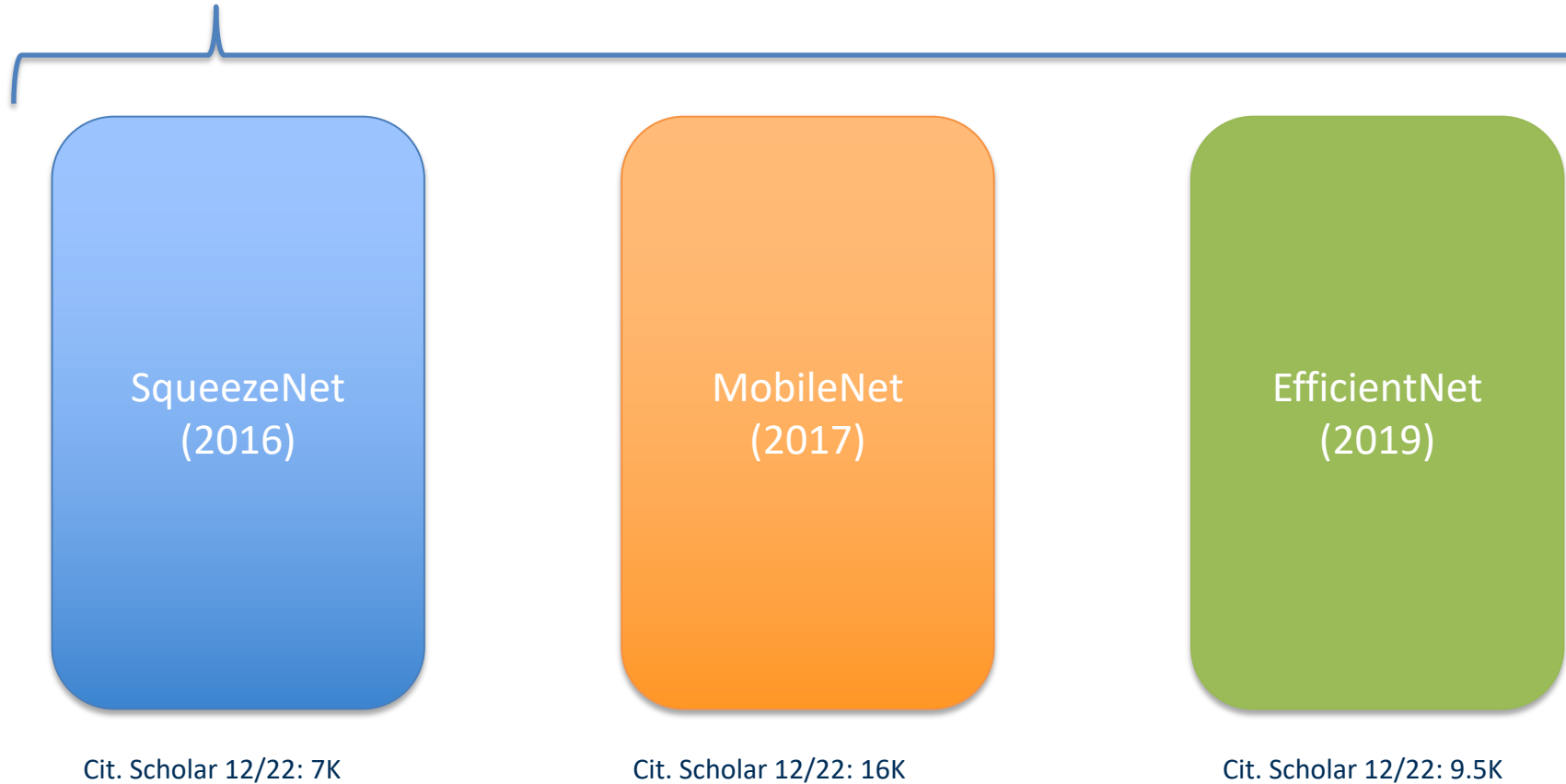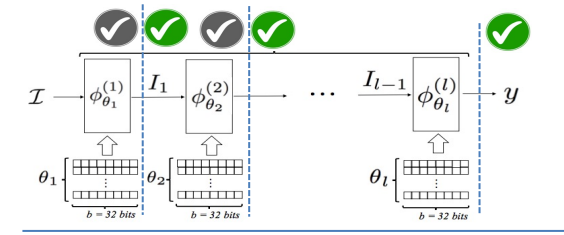  - ✓ **Architectures and EECs**
  - ✓ **Learning EENNs**

# 5) Deep Learning for Embedded and Edge AI (8h)

Adaptive mechanisms for Embedded and Edge AI

# 6) From the unit to the ecosystem perspective (4h)

# 7) The ethical perspective (2h)



"Ethics of Design and Values: Solutions and Trade-offs in H-IoT and Beyond"
**Prof. Viola Schiaffionati – Prof. Stefano Canali**

# The labs (20h)

1. Intro/review on Embedded Systems
2. Intro/review on Deep Learning with Tensorflow

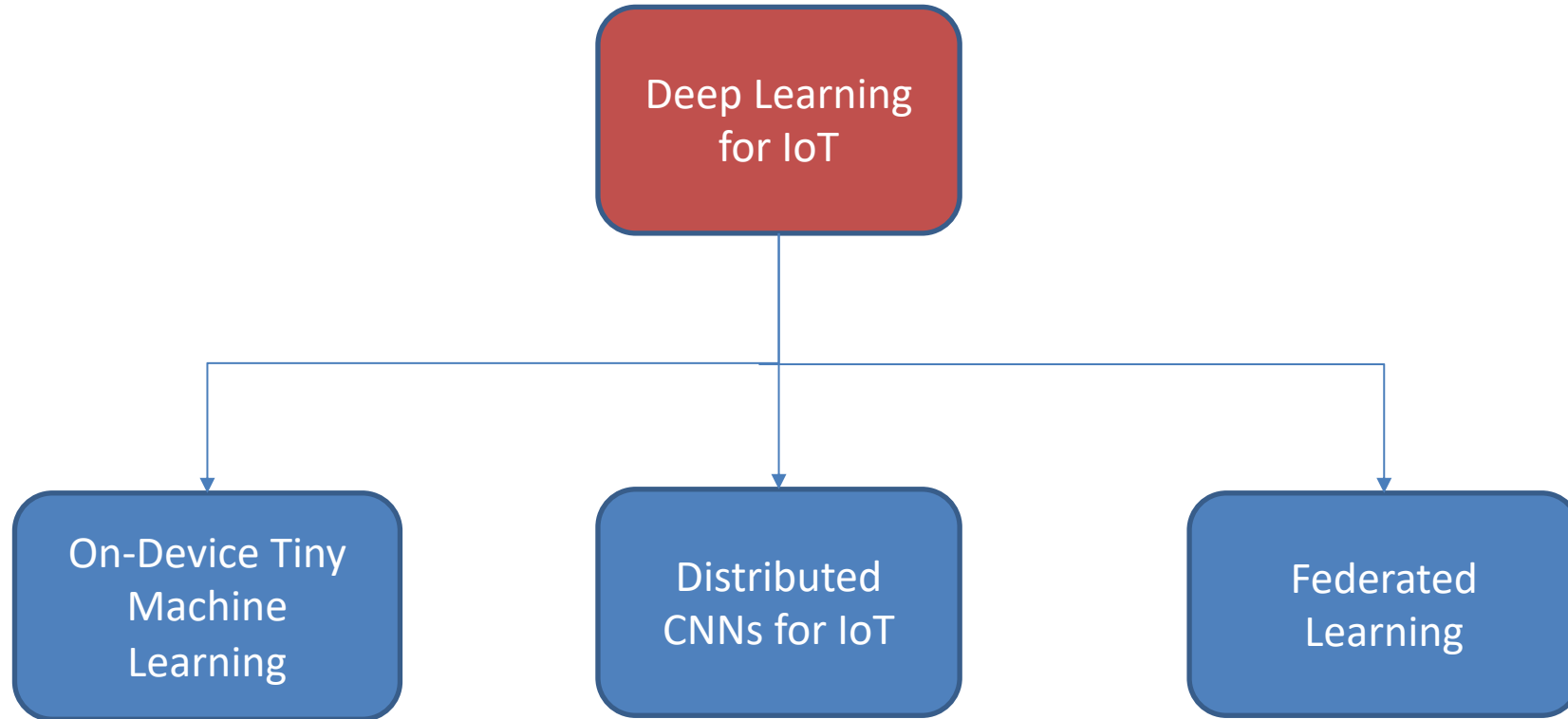*Introduction and review of needed background*

3. TF lite optimizations: quantization and pruning
4. Training Keyword Spotting – Microphone
5. Deploying Keyword Spotting – Microphone
6. Training and Deploying Visual Wake Word – Camera

*Core TinyML Lectures*

7. Data Collection and Engineering with Edge Impulse – Camera
8. Options for deployment: TFLM and Edge Impulse
9. Training and Deploying Anomaly detection – Accelerometer

*Innovative Tools and topics for the future of TinyML*

Tools employed: Google Colab, Edge Impulse, Arduino IDE, TFLM

# TinyML kit – The Arduino Nano 33 BLE sense



Legend:
- Ground
- Power
- LED
- Internal Pin
- SWD Pin
- Digital Pin
- Analog Pin
- Other Pin
- Microcontroller's Port
- Default

RGB LED

| BUILT_IN LED | P0.13 | LED_BUILTIN |
| Power | P1.09 | LED_PWR |

| SCK | P0.13 | D13~ | | ~D12 | P1.08 | CIPO |
| | | +3V3 | | ~D11 | P1.01 | COPI |
| | | AREF | | ~D10 | P1.02 | |
| | P0.04 | A0 | | ~D9 | P0.27 | |
| | P0.05 | A1 | | ~D8 | P0.21 | |
| | P0.30 | A2 | | ~D7 | P0.23 | |
| | P0.29 | A3 | | ~D6 | P1.14 | |
| SDA | P0.31 | A4 | | ~D5 | P1.13 | |
| SCL | P0.02 | A5 | | ~D4 | P1.15 | |
| | P0.28 | A6 | | ~D3 | P1.12 | |
| | P0.03 | A7 | | ~D2 | P1.11 | |
| | | +5V | | GND | | |
| | | RESET | | RESET | | |
| | | GND | | RX | P1.10 | |
| | | VIN | | TX | P1.03 | |

# Exam

- The exam will consist in **<u>two parts</u>**:
    1. **Written exam** (16 points) comprising questions (closed/open) about the topics of the course
    2. **Project** (16 points):
        - Your own idea with our own hardware
        - Max 2 people
        - Delivered at the exam dates
        - Code + presentation
        - Evaluation will take into account:
            - The "market" perspective (5 points)
            - The "technological" perspective (6 points)
            - The "ethical" perspective (5 points)

# Selected projects of the course

# Challenges and opportunities

- Heterogeneity of the students backgrounds
- Fast evolution of the technology
- Keep the correct trade-off between ambition and implementability in the students' projects

- Strong connection between research activity and teaching
- The presence of a "physical" lab to carry out the projects
- Combining theory with implementation
- Strong technical aspects with ethical flavor