# EMBEDDED AI: INTELLIGENCE ON DEVICES

Serge Pacome Bosson, Ing. Radek Holota, Ph.D, prof. Ing. Václav Šmídl, Ph.D.
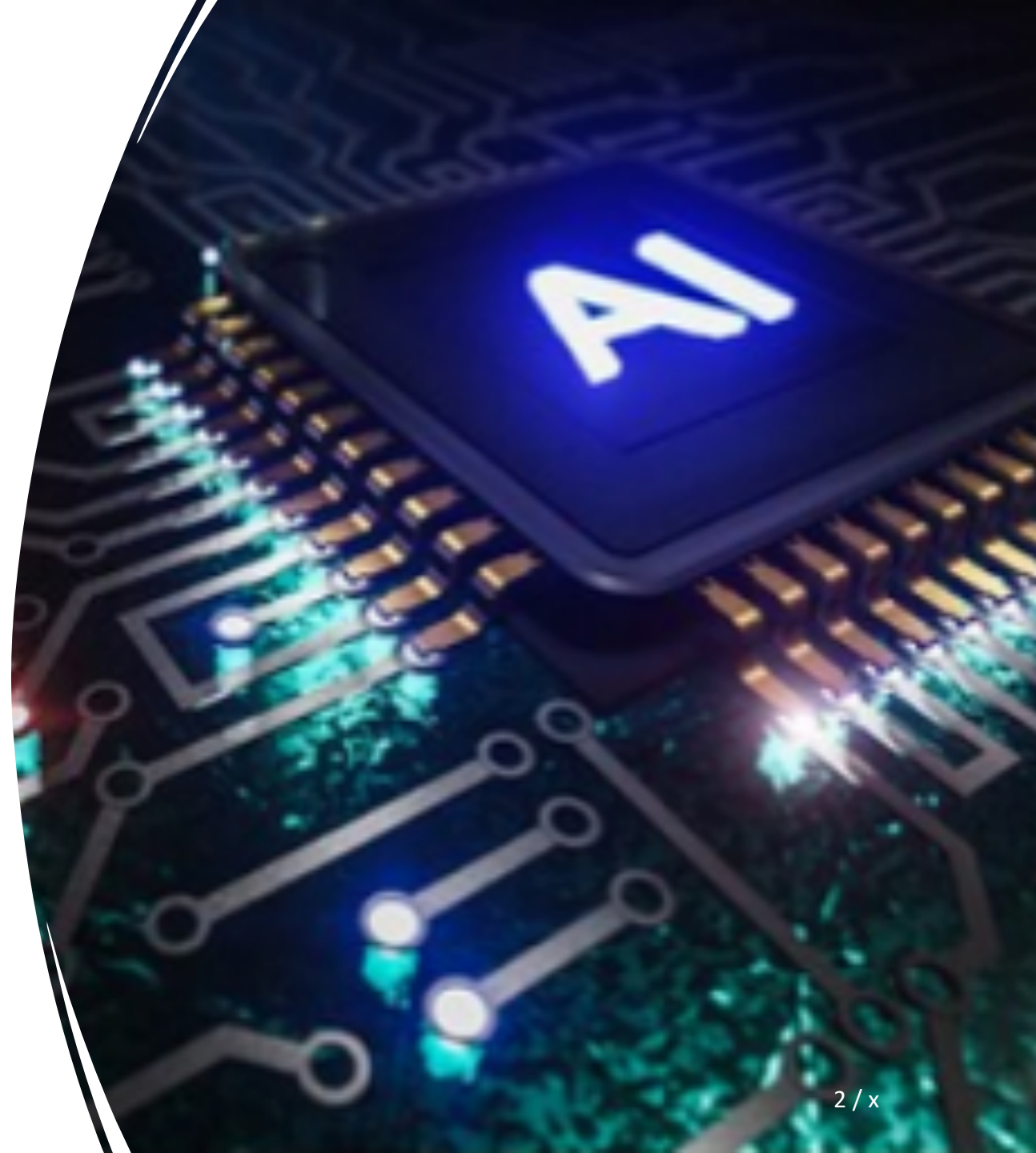
# MOTIVATION

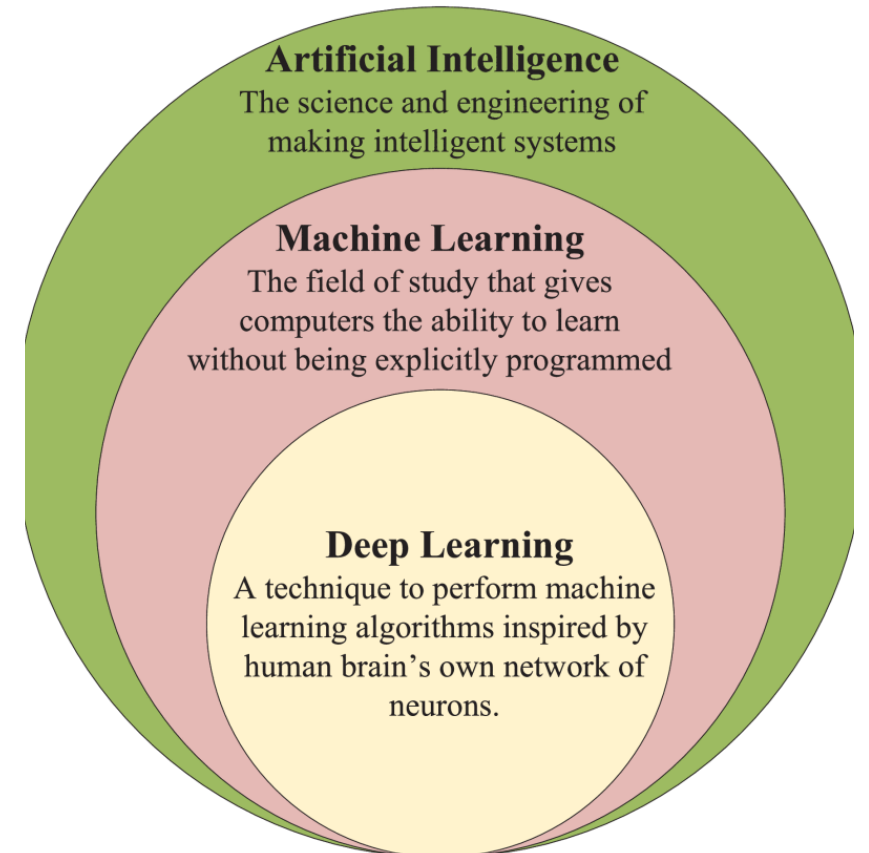▌ Deploy Machine Learning Algorithms onto an Embedded System.
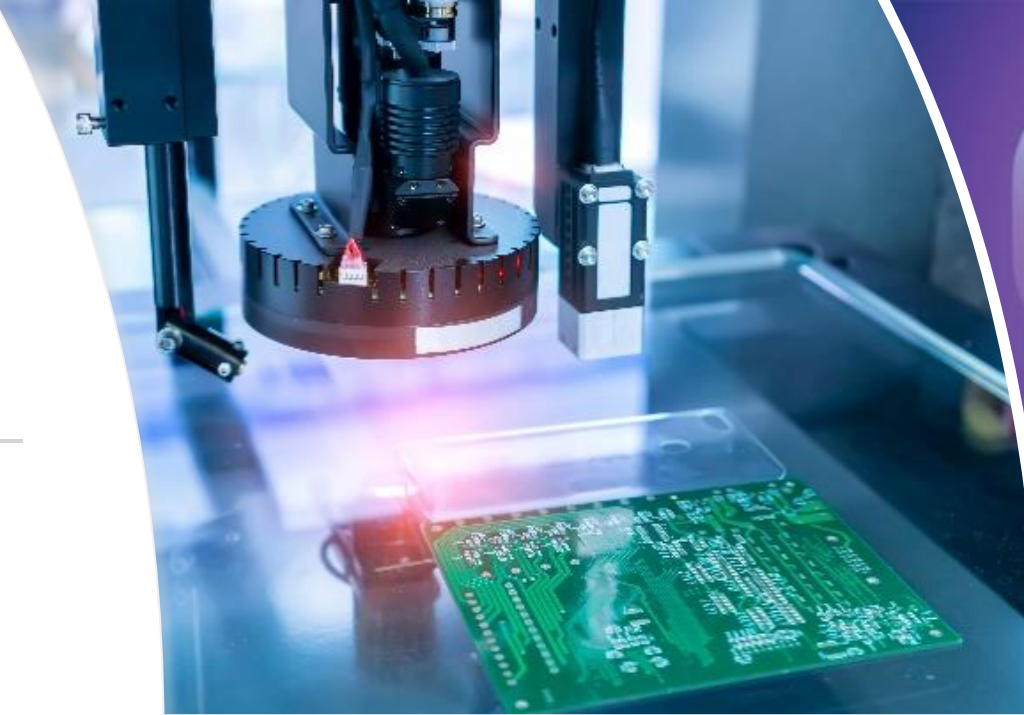
# WHAT IS AI AND DEEP LEARNING?

▌ Artificial intelligence (AI): Using human-like intelligence to solve tasks.

▌ Machine learning (ML): Algorithm uses data to find patterns.

▌ Deep learning (DL): Very large algorithms using raw data input.

▌ Machine learning provides significant advantages over classical computing:
   - Scalability
   - Less R&D effort
   - More accurate

**Artificial Intelligence**
The science and engineering of making intelligent systems

**Machine Learning**
The field of study that gives computers the ability to learn without being explicitly programmed

**Deep Learning**
A technique to perform machine learning algorithms inspired by human brain's own network of neurons.

# WHERE DO WE USE AI ?

- Drive assistant

- Machine Vision.

- Fault diagnosis.

- Robotics.

- Security and
  home Automation cameras.

- Speech recognition,
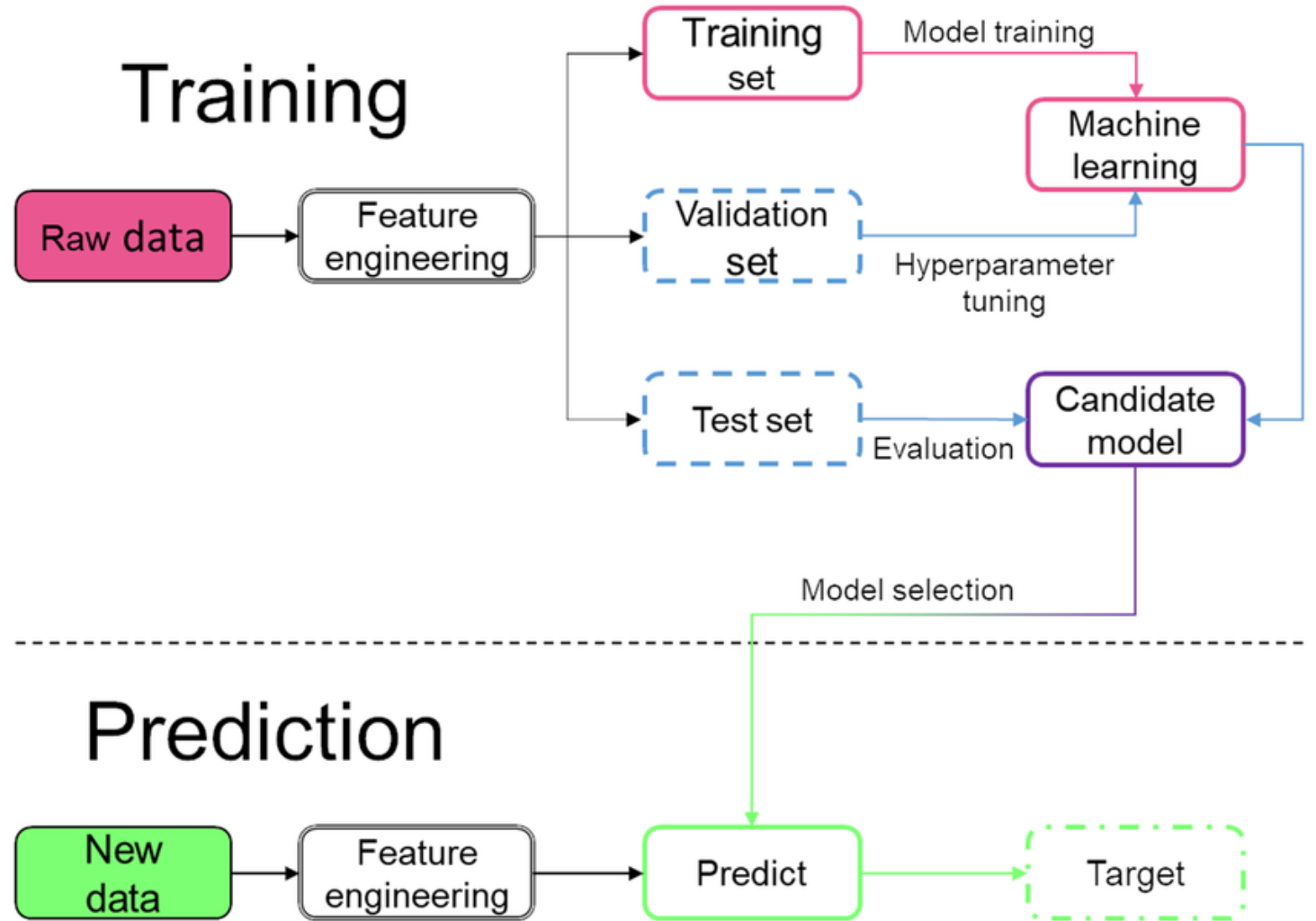  text analysis, translation.
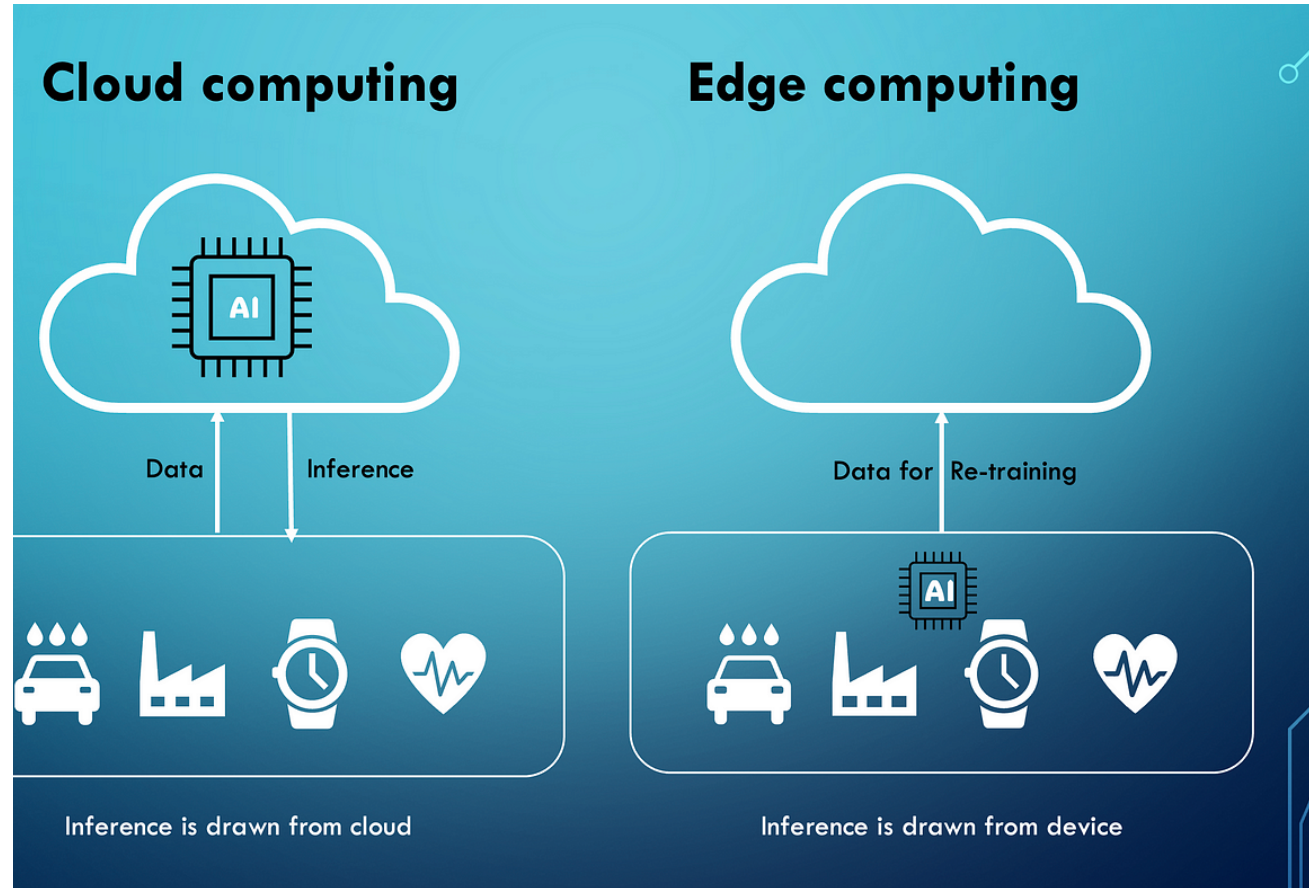
- And Many More!

# PRINCIPLE

- Collect data.
- Design a ML model.
- Train the model.
- Move trained model to platform for inference.

# AI INFERENCE

▌ Edge AI: run ML models where data is generated.
  - Algorithms run on embedded systems.

▌ Cloud AI: run ML models on cloud servers.
  - Algorithms run on data centers.

▌ Benefits of Edge AI:
  - Reduced latency.
  - Improved privacy and security.
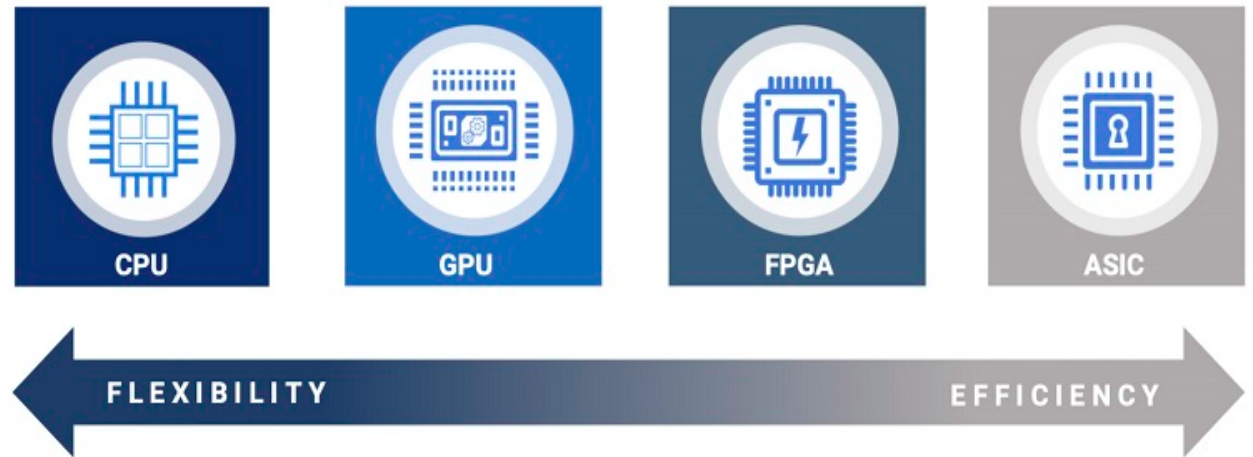  - Enhanced energy efficiency.
  - Real-time decision-making.



**Cloud computing**

Data · Inference

Inference is drawn from cloud

**Edge computing**

Data for Re-training

Inference is drawn from device

# HARDWARE FOR EMBEDDED AI

The brain of an embedded AI device is usually based on:

▌ Microprocessor (CPU).

▌ Graphical processing unit (GPU).

▌ Field programmable gate array (FPGA).

▌ Application specific integrated circuit (ASIC).



CPU    GPU    FPGA    ASIC

FLEXIBILITY ← → EFFICIENCY

# AI INFERENCE



**Embedded AI | Application development flow**

AI model Training

1. Data preparation — Time consuming
2. Model Designing — AI expertise
3. Model training — Data scientist

ML Experts · Data Scientists

AI inferencing

4. Model Export — Vendor specific
5. Model Deployment — Vendor specific
6. Application integration — Vendor specific tools, frameworks & APIs

Embedded Developers

# ML COMPILER FRAMEWORK FOR CPU AND GPU: OSRT TECHNOLOGY

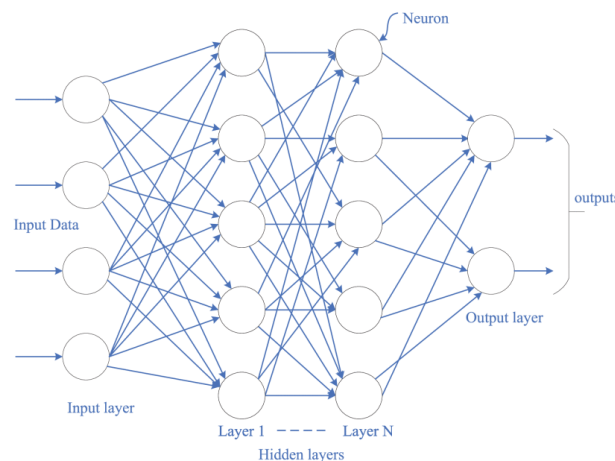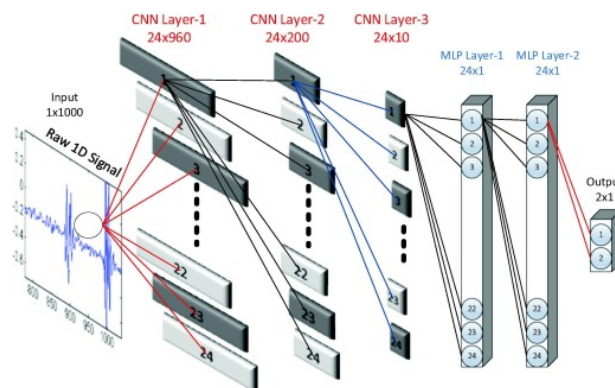# CAS STUDY: FUNCTION APPROXIMATION IN EMBEDDED SYSTEMS

Function approximation based on machine learning algorithms can find practical application in electrical engineering:

- Flux linkages approximation.
- Online fault diagnosis.
- AI robotics.
- Advanced control in power electronics and drives.
- Optimization problems.
- Etc…

# CAS STUDY: FUNCTION APPROXIMATION IN EMBEDDED SYSTEMS

Approximators that admit efficient implementation in conventional industrial computer systems:

- Multilayer perceptron (MLP).
- 1D Convolutional neural network (1D CNN).
- Piece-wise affine (PWA).
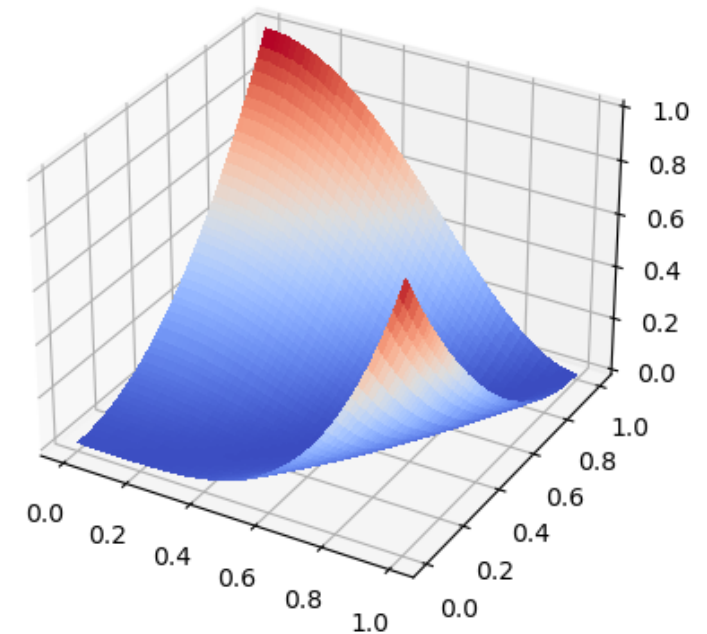- Lattice interpolated look-up table (LUT).

# DATA PREPARATION

Evaluate the regression power of the reviewed approximators on the classical optimization problem called Rosenbrock's valey:

- Generated 10K points from the uniform distribution [0 1]^D.
  - D: Number of dimension: 2D 5D 8D 15D.
- Compute D-dimensional Rosenbrock's function and scale them in the range of [0 1].



2-D Rosenbrock's valley function

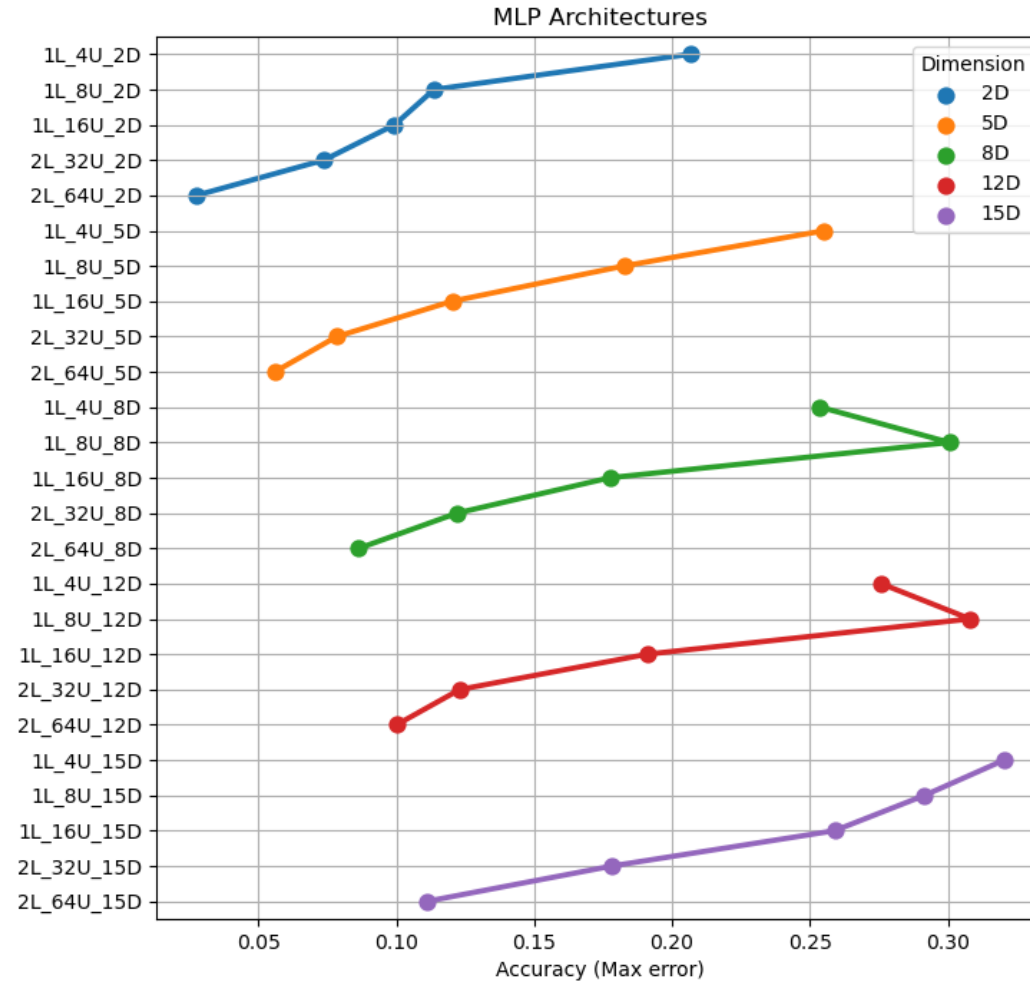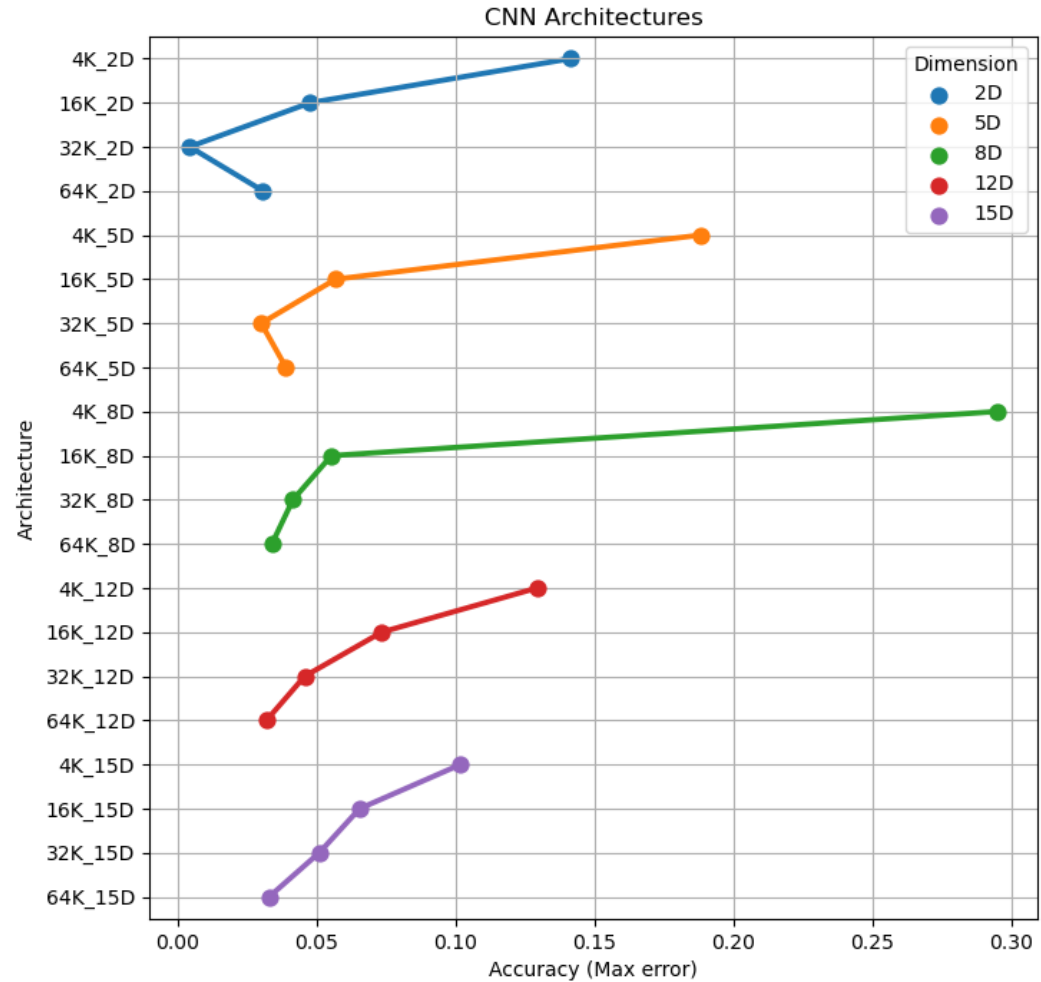# MODEL DESIGN

Multi-dimensional ML model design:

- MLP design using Tensorflow and Keras library:
    - Varying depth (number of layer) and width (number of neuron in each layer).
- 1-D CNN design using Tensorflow and Keras Library:
    - 1 convolution layer and variations of number of filter with fixed kernel size equal to 2.
- Lattice LUT design using Tensorflow lattice library:
    - Varying the lattice sizes.
- PWA design using Python code source available at https://github.com/bemporad/PyPARC.git and modify the code under Apache 2.0 license to add metrics for accuracy comparison with other ML models:
    - Varying number of partitions.

# AI MODEL TRAINING: MLP Architectures



MLP Architectures

# AI MODEL TRAINING: CNN Architectures



CNN Architectures
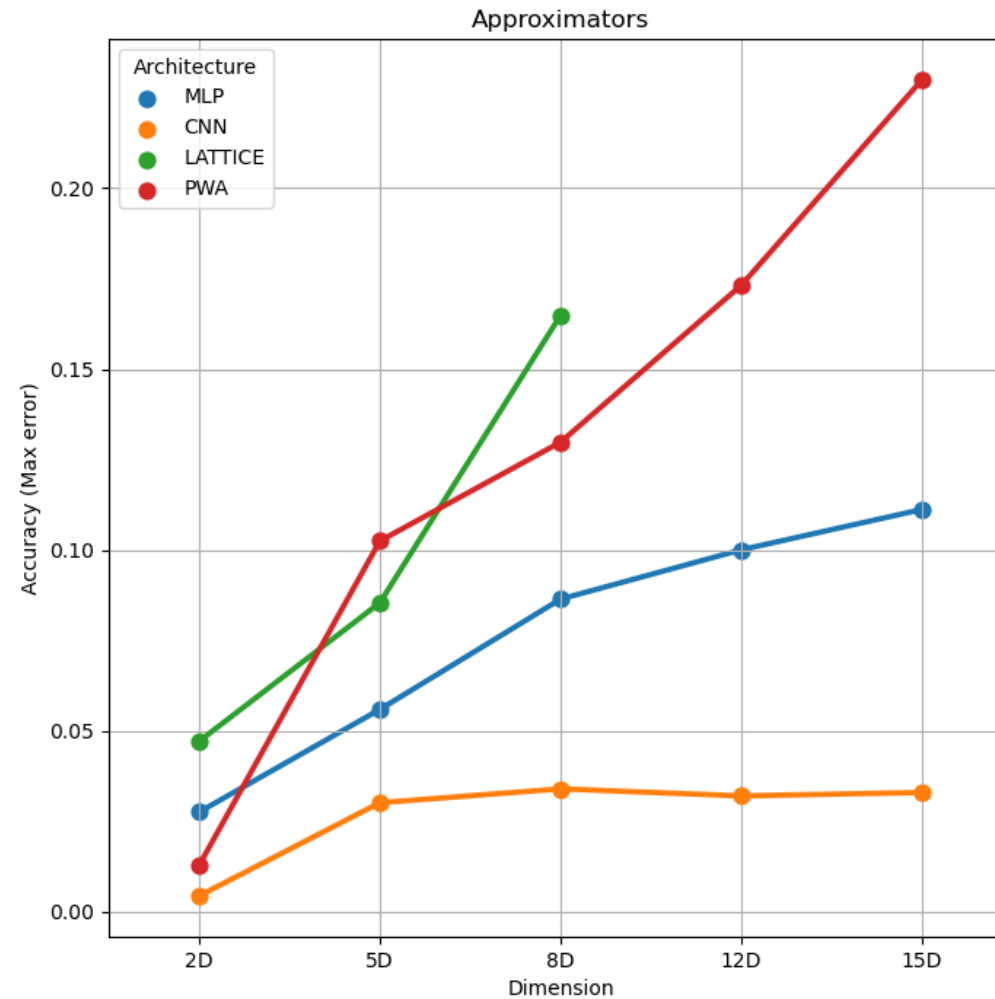
# AI MODEL TRAINING: PWA Architectures



PWA Architectures

# AI MODEL TRAINING: LATTICE Architectures



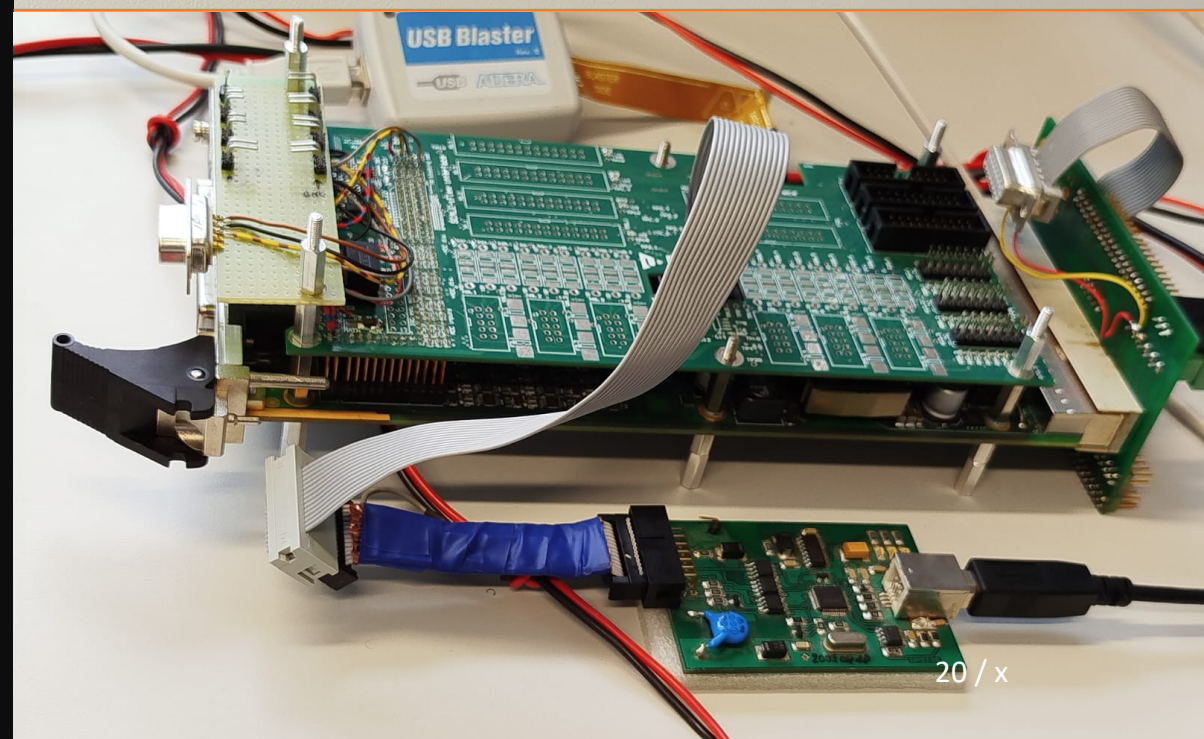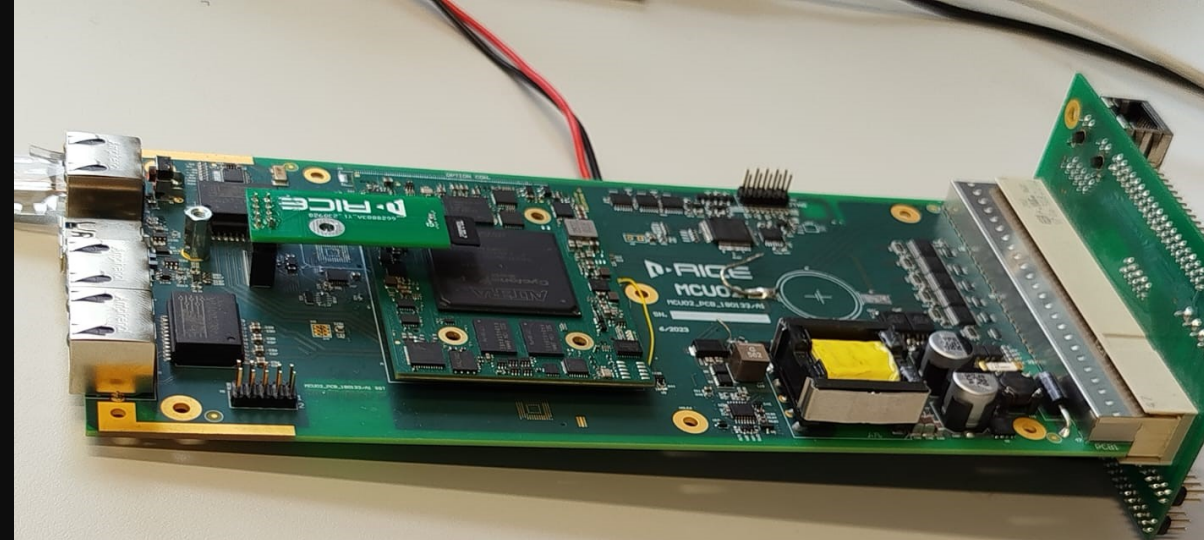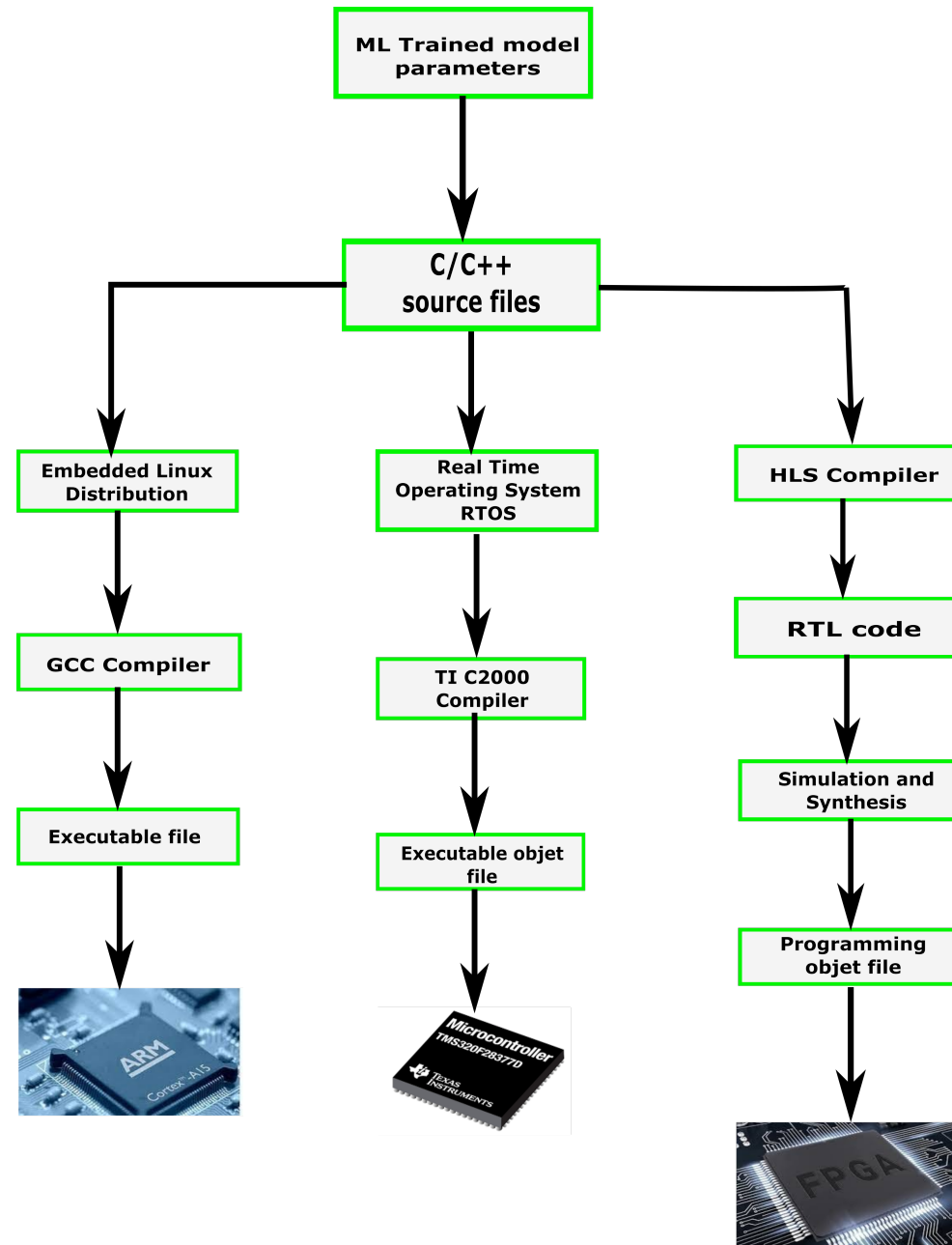LATTICE Architectures

# AI MODEL TRAINING: APPROXIMATORS

# AI MODEL INFERENCE

- Our industrial computer systems are based on:
  - SoC FPGA Cortex-A Arm – embedded Linux.
  - Cyclone V FPGA.
  - Texas instruments dual core MCU.
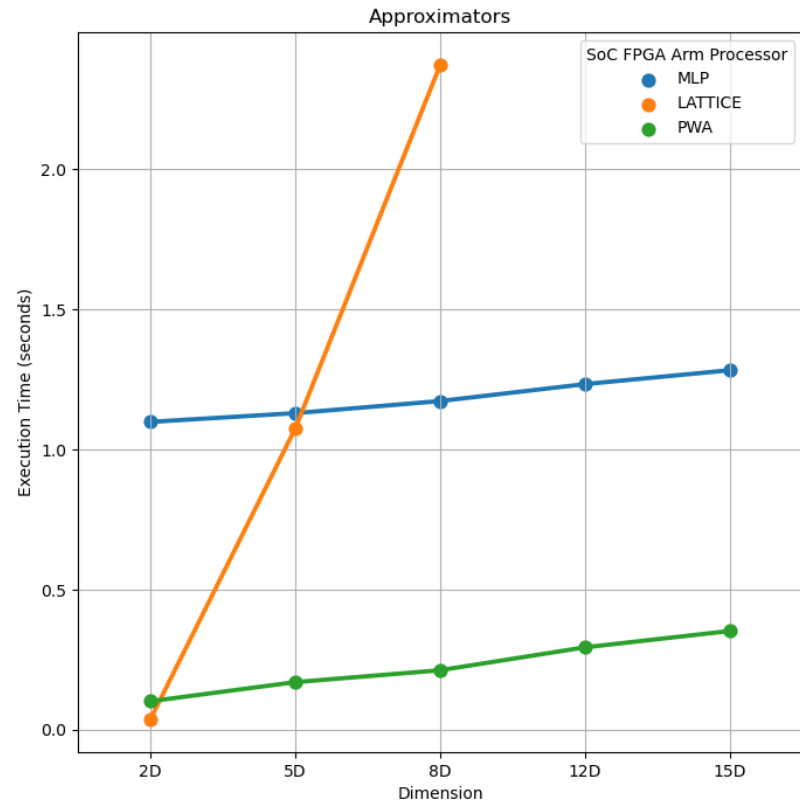
# APPROXIMATORS IMPLEMENTATION: SOFTWARE DEVELOPMENT FLOW

ML Trained model parameters

C/C++ source files

Embedded Linux Distribution

Real Time Operating System RTOS

HLS Compiler

GCC Compiler

TI C2000 Compiler

RTL code

Executable file

Executable objet file

Simulation and Synthesis

Programming objet file

# AI INFERENCE: SoC FPGA

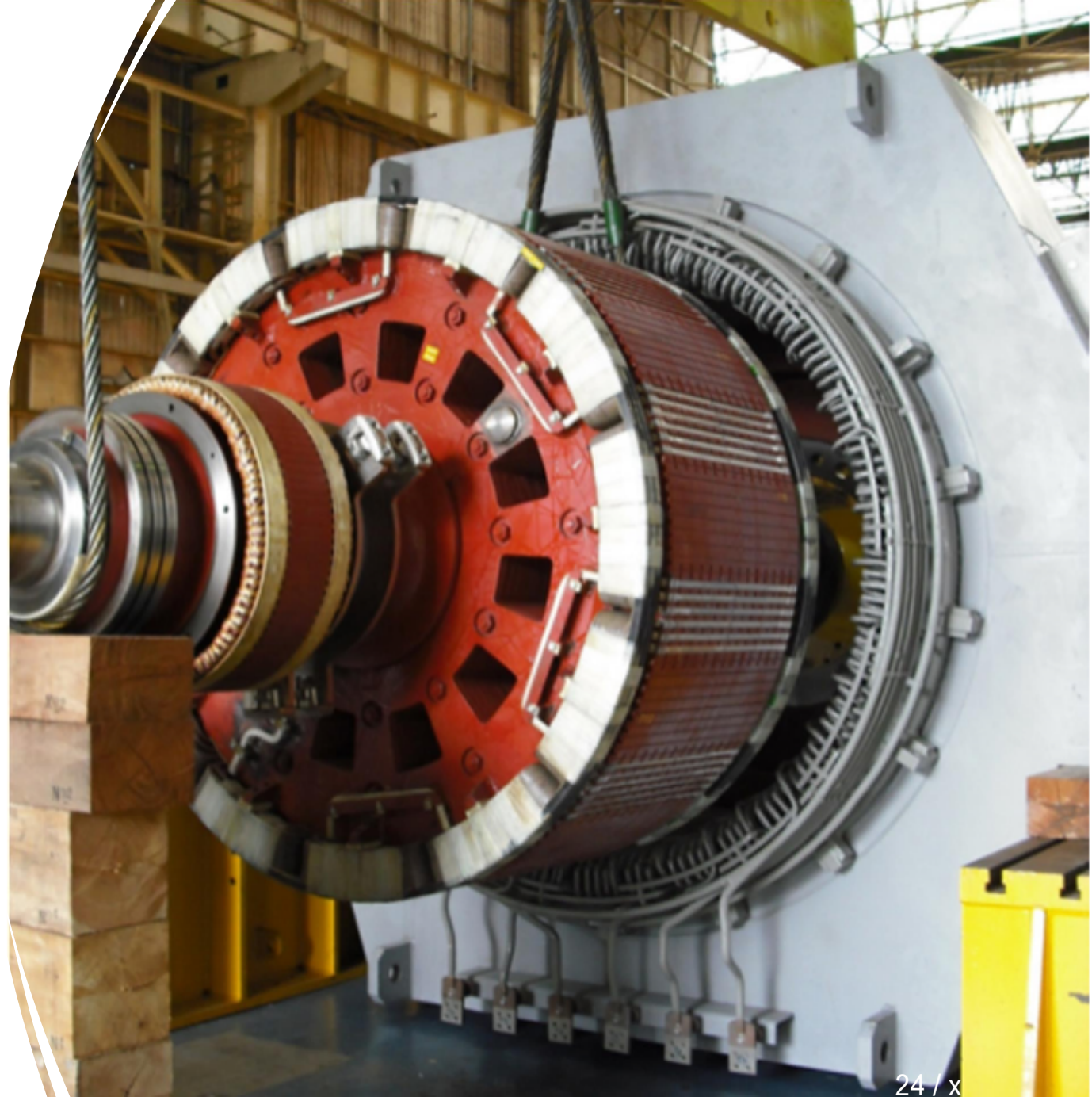# AI INFERENCE-SoC FPGA: Execution time vs Dimension: 4K streaming data.

# APPLICATION EXAMPLE: DIAGNOSTIC OF SYNCHRONOUS GENERATOR

Main project is online diagnostic of large SG with power range of hundred MW used in a power plant:

- Collect data from sensors installed inside the SG.

- Design a binary classification for fault detection or design a regression model for fault prediction/state of the health estimation.

- Perform some feature engineering on the raw data then use a MLP to design the diagnostic model.

- Or use a 1DCNN directly on the raw data to design a diagnostic model.

- Deploy the model inference on embedded systems.

# Thank You for your attention!

**Serge Pacome Bosson**
377 634 408
bosson@fel.zcu.cz
fel.zcu.cz