

Introduction to Edge Al

Ireland AI Roadshow 2025

George Dickey

Global Technology Leader

WUVN

George Dickey

Global Technology Leader High-end Processing & Al

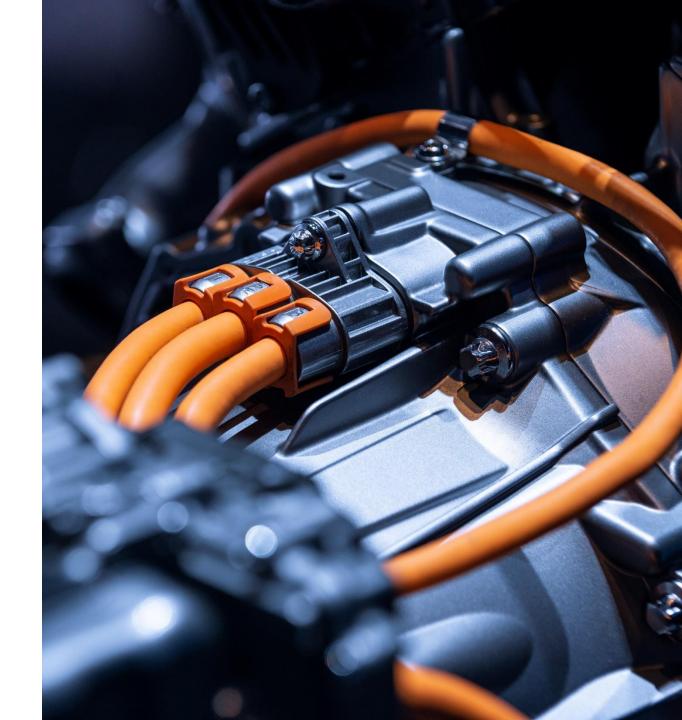
- Deep technical background
 - Digital ASIC design
 - FPGA
 - Embedded Software & DSP
 - Artificial Intelligence
- Across multiple roles
 - Designer
 - Engineering Manager
 - CTO
- In different scaled organisations
 - Independent Consultant
 - SME
 - Multinational



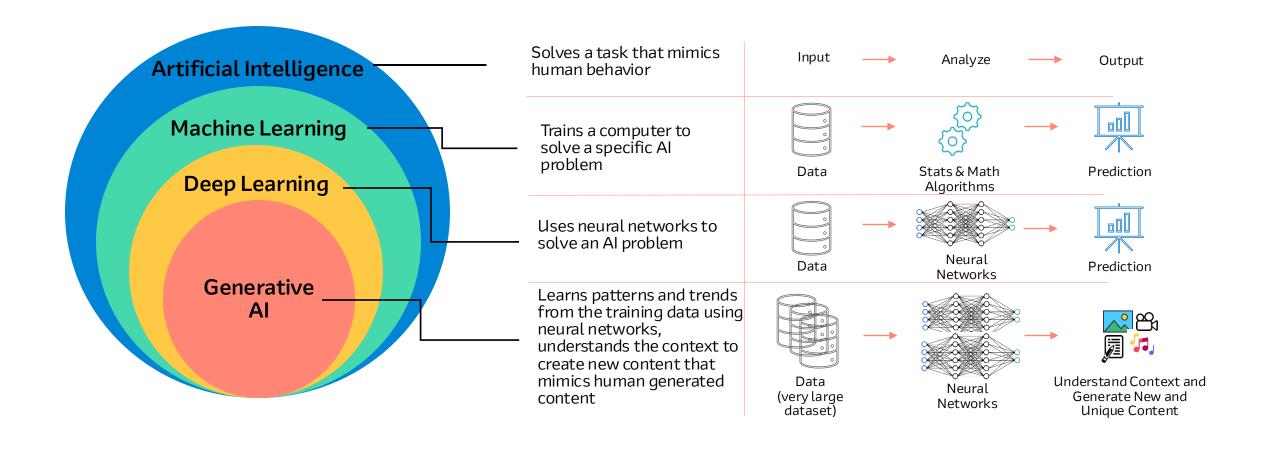
Agenda

Artifical Intelligence

- What is AI and why use it?
- Why at the Edge?
- Device Selection
- Final thoughts



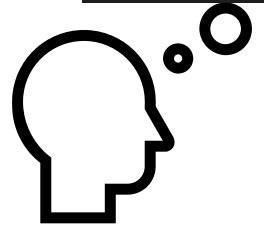
What is an AI model?



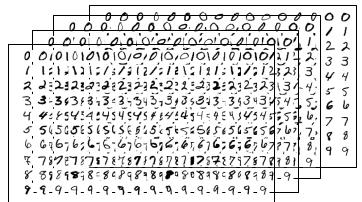
Edge AI is primarily Deep Learning

Why use AI? 000000000000000 1///// 2222222222 3333356332 9 4 4 4 5 5 5 5 5 5 5 5 5 5 5 66666666666666 7777777777777 888888888888

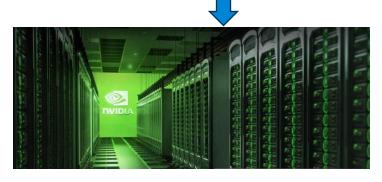
Design Engineer



Deep Learning



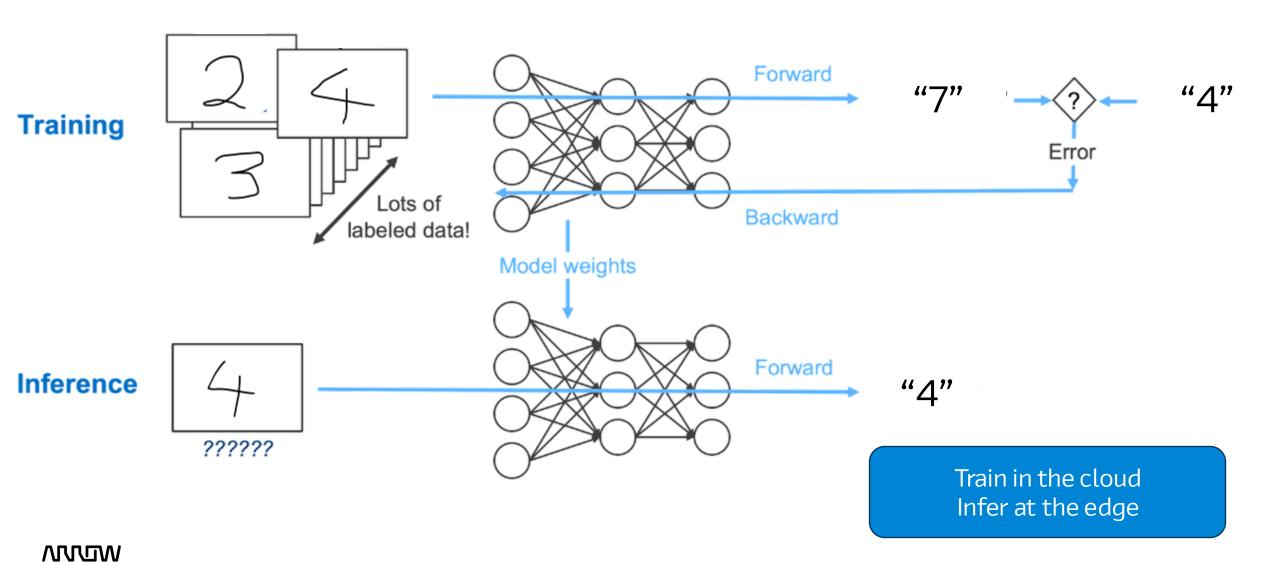
Data



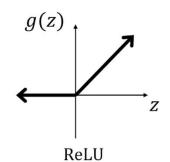
Compute

Deep Learning solves complex problems by using lots of data and compute resources rather than human derived algorithms

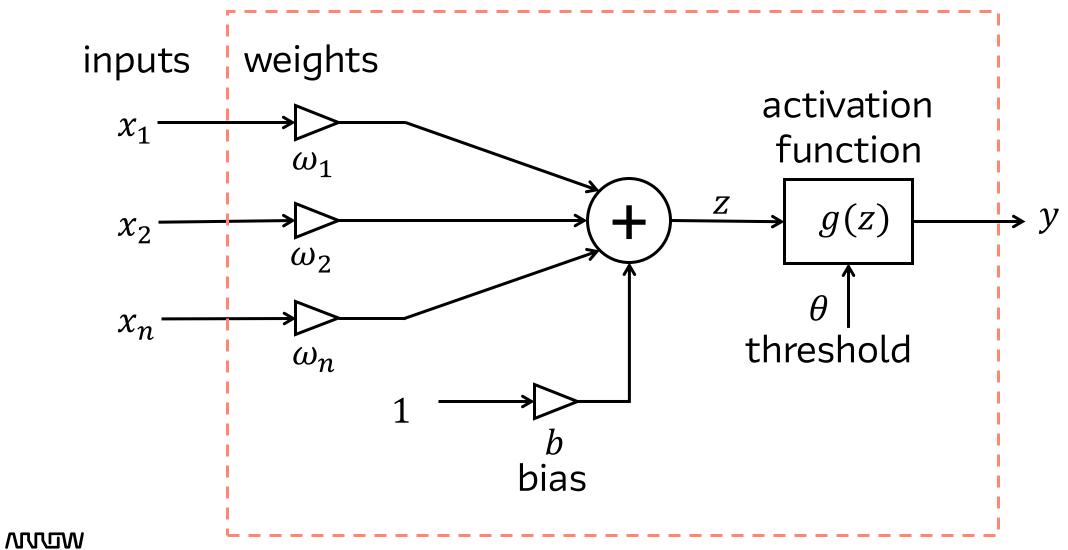
Training and Inference



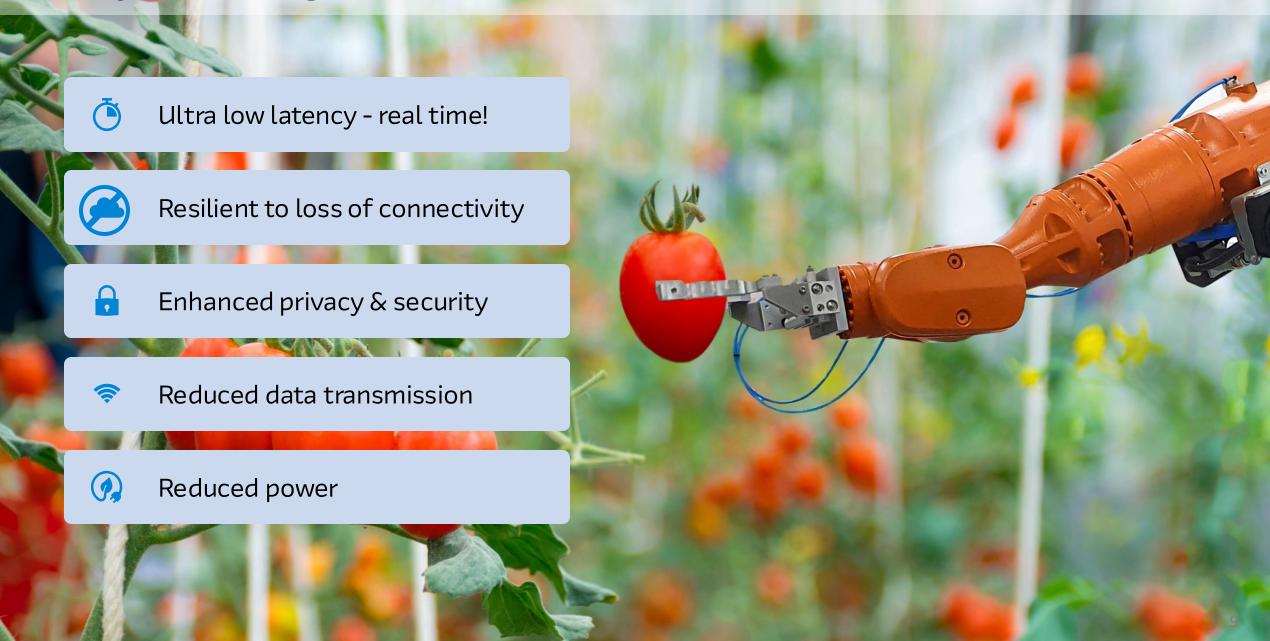
It's all just maths



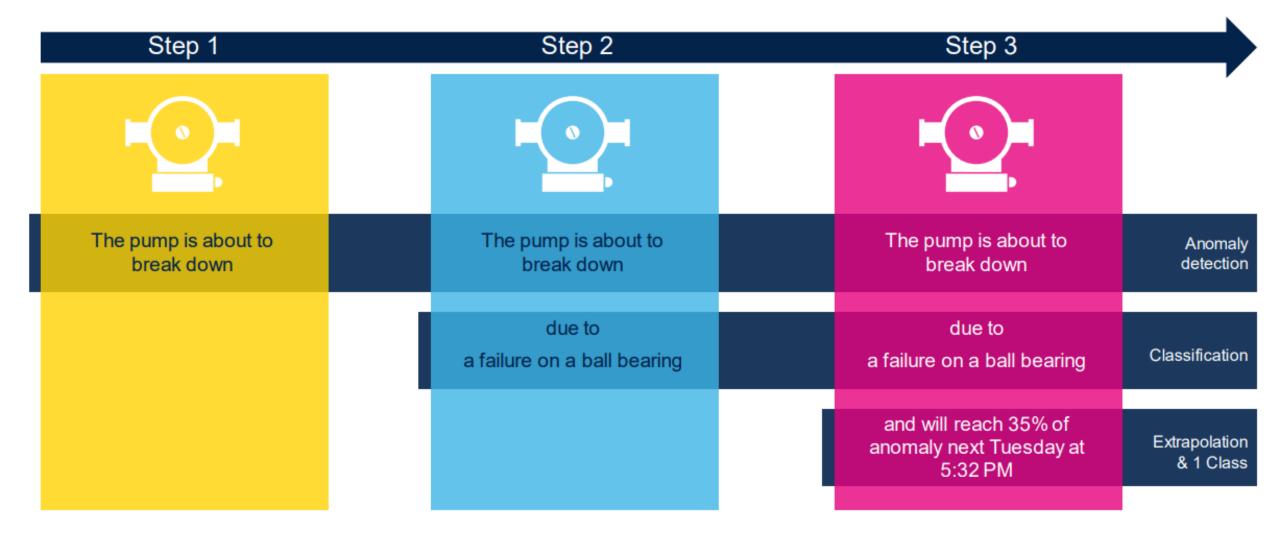




Why the Edge?



Predictive Maintenance





AI/ML development Flow

The Five steps to realize AI

Data Engineering Model Training Optimisation Deployment Monitoring

Monitoring

Specify, collect, clean and prepare data for model training and testing

Selecting, training and testing the model

Optimise the model for the edge, to improve performance and reduce footprint

Deploying and running the model on HW platform

Monitoring the model for drift, basis and accuracy

Dataset Collection

Open problem



- Environment, factors beyond control
- Unknown, various situations or objects
- Difficult tail cases

Closed problem



- Controlled environment
- Well-defined, known objects
- Easier to address tail cases

Data
Engineering

Model
Training

Model
Optimisation

Deployment

Monitoring

Transfer Learning

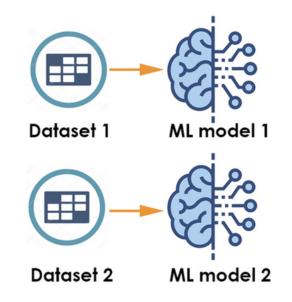
Why

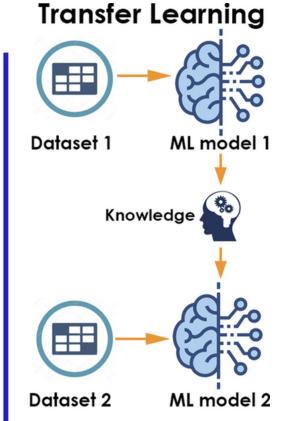
- Reduced Training Time
- Improved Performance
- Data Efficiency

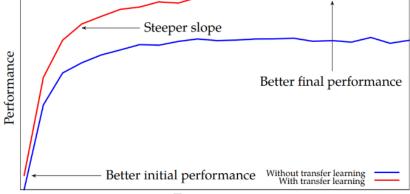
How

- Start with a pre-trained model (baseline)
- Lock all but the last few layers (or add additional layers)
- Re-train the unlocked layers of the mode
- Unlock and fine-tune the whole model

Training from scratch







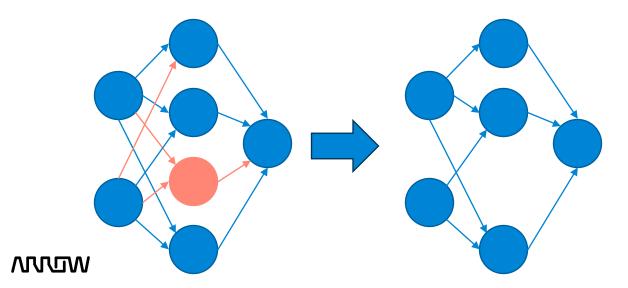


Quantization

- Typically floating point (FP32) to integer (Int8)
- Int4, ternary and even binary are being researched for further reduction in memory and complexity

Pruning

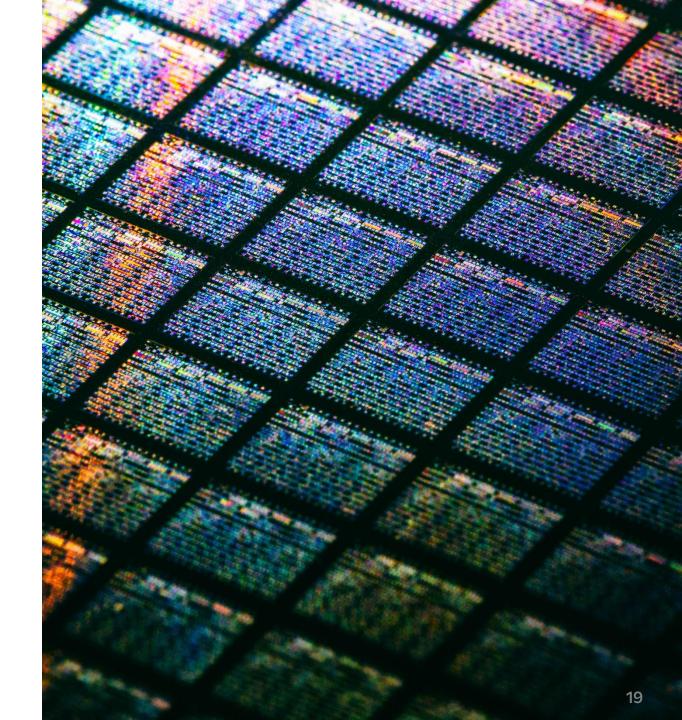
- Connections with weights that are close to 0
- Nodes whose output is close to 0

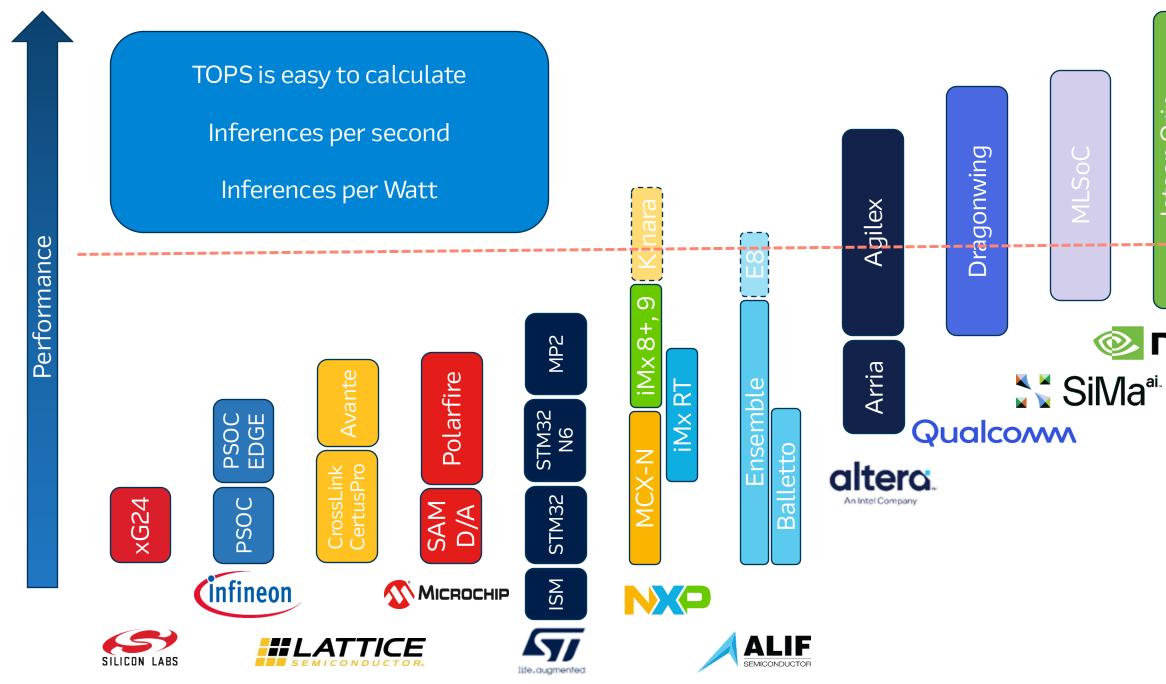


Device selection

Factors

- Performance GOPS / TOPS
- Through put
- Latency
- Al acceleration
- Development Tools
- Application specific cores
- I/O
- Power
- Cost
- Community







etson

NVIDIA.

Neural Processing Unit (NPU)

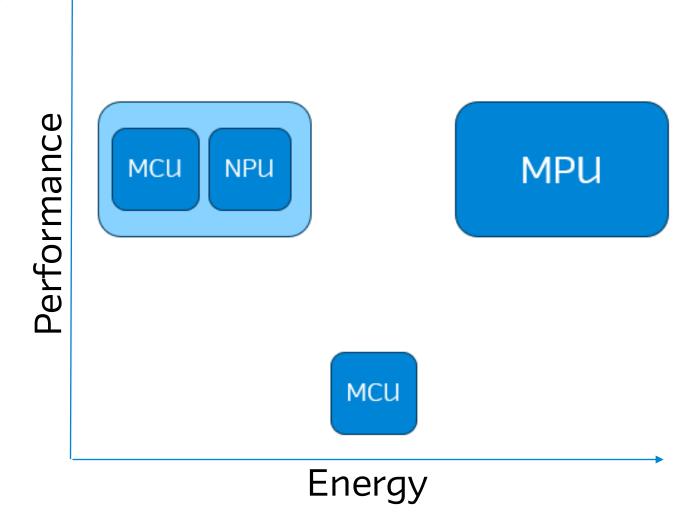
An NPU is dedicated hardware that is very efficient at AI operations

For AI tasks an NPU

- Provides higher performance
- For less energy

Microcontroller with NPU vs Microprocessor

- Similar AI performance
- Less power
- Less memory
- Lower total BOM cost



Neural Processing Unit (NPU)

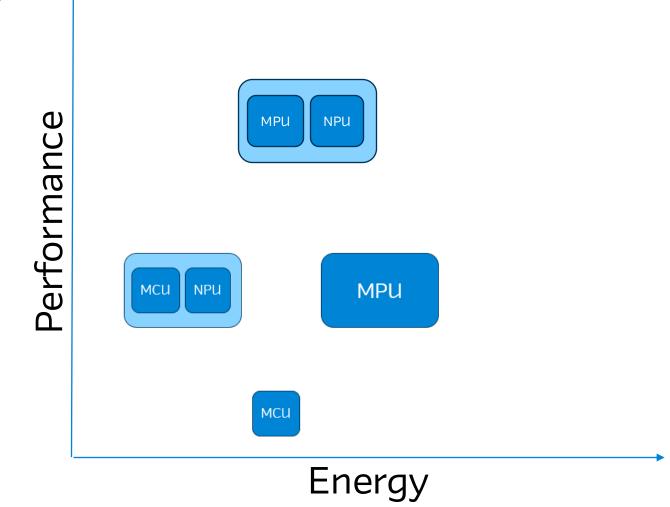
An NPU is dedicated hardware that is very efficient at AI operations

For AI tasks an NPU

- Provides higher performance
- For less energy

Microprocessor with NPU vs Microprocessor

- Higher Al performance
- Less power
- Less memory
- Lower total BOM cost



Neural Processing Unit (NPU)

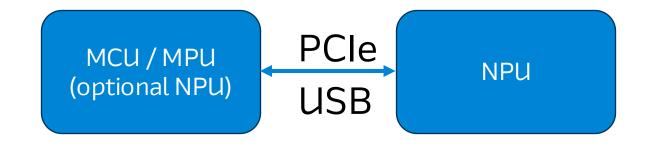
An NPU is dedicated hardware that is very efficient at AI operations

For AI tasks an NPU

- Provides higher performance
- For less energy

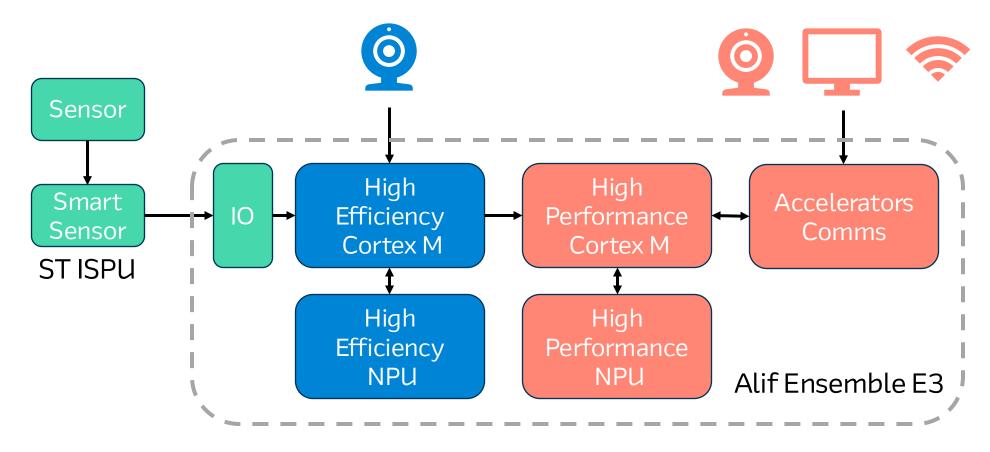
MCU / MPU with external NPU

- Substanially higher AI performance
- Low power per inference



Hardware and multiple models

It's all about sleeping when you can



Always On

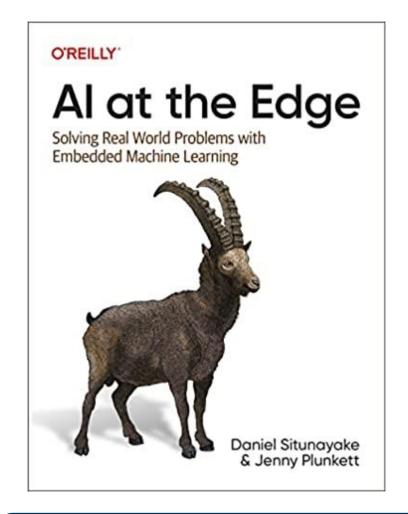
High Efficiency

High Performance

Final thoughts

For your embedded Al journey

- Have a go! <u>Edge Impulse</u>
- *Transfer Learning* Make use of pre-trained models
- Embedded Engineers
 - Use an end-to-end software studio if possible
- AI Specialists
 - Al tools for data engineering and model training
 - Supplier tools for optimization and deployment
- Use existing hardware (devkits, SOMs) where you can, especially for prototyping.
- Recommended Book <u>Al at the Edge</u>, O'Reilly



We can help!

Arrow, driving innovation forward!



arrow.com/edge-ai

arrow.com/edge-ai-webinars



linkedin.com/in/gdickey