**Optimized Neural Networks on STM32** with STM32Cube.AI
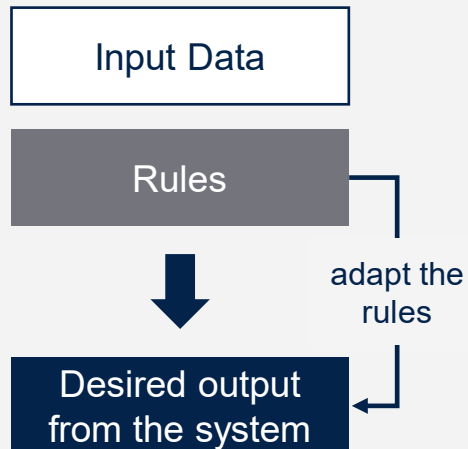
STM32 Cube.AI

# Introduction to Edge AI

# A new way to add environment awareness to your products
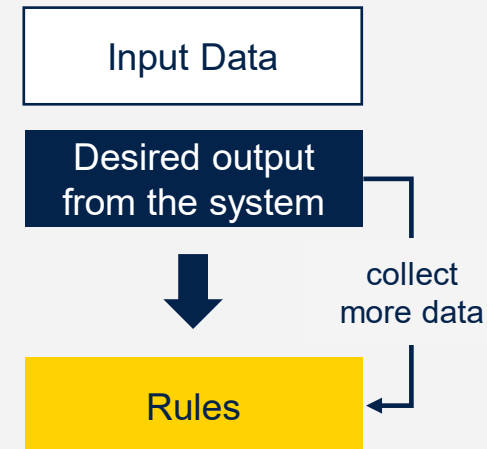
## From rule-based engineering to data-driven engineering

### Standard programming
**Handcrafted rules based on experience**

```
Input Data
   │
Rules ──────┐
   │    adapt the
   ▼      rules
Desired output ◄──┘
from the system
```

- Requires domain expertise to code
- Need to rewrite if environment evolves

### Machine Learning
**Rules learnt from real-world data**

```
Input Data
   │
Desired output ──────┐
from the system  collect
   │            more data
   ▼
Rules ◄──────────────┘
```

- Generate code from real-world observations
- Re-learn from data if environment evolves

# Distributed Artificial Intelligence approach

**Leverage billions of devices at the Edge!**



| Thousands | **Data Center Cloud** Analytics, storage, compute |
| Millions | **Edge Nodes** IoT gateways, micro datacenters |
| 100 Billions | **Edge Things** Real time, local processing |

Ultra-low-power devices and sensors

STM32

# Artificial intelligence at the Edge

## Moving part of Artificial Intelligence closer to the data acquisition brings several benefits
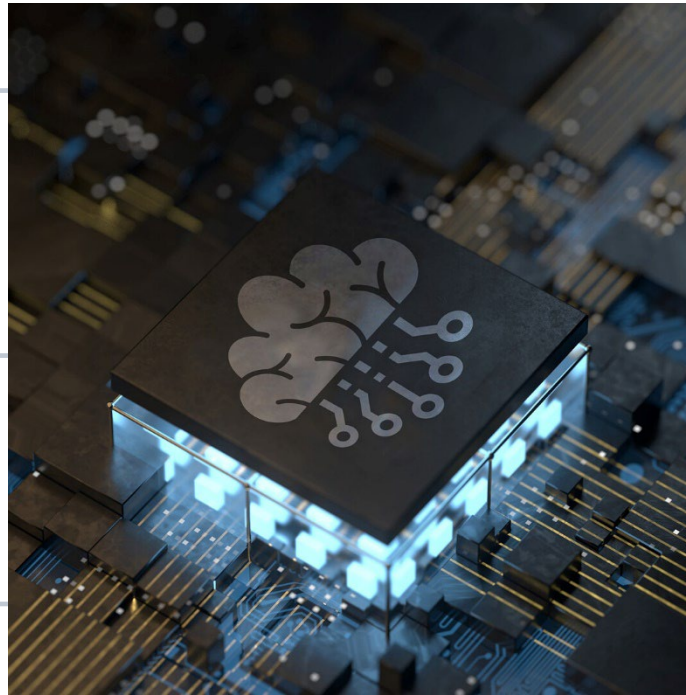
**Ultra-low latency**
Real-time applications

**More reliability**

**Security of data**
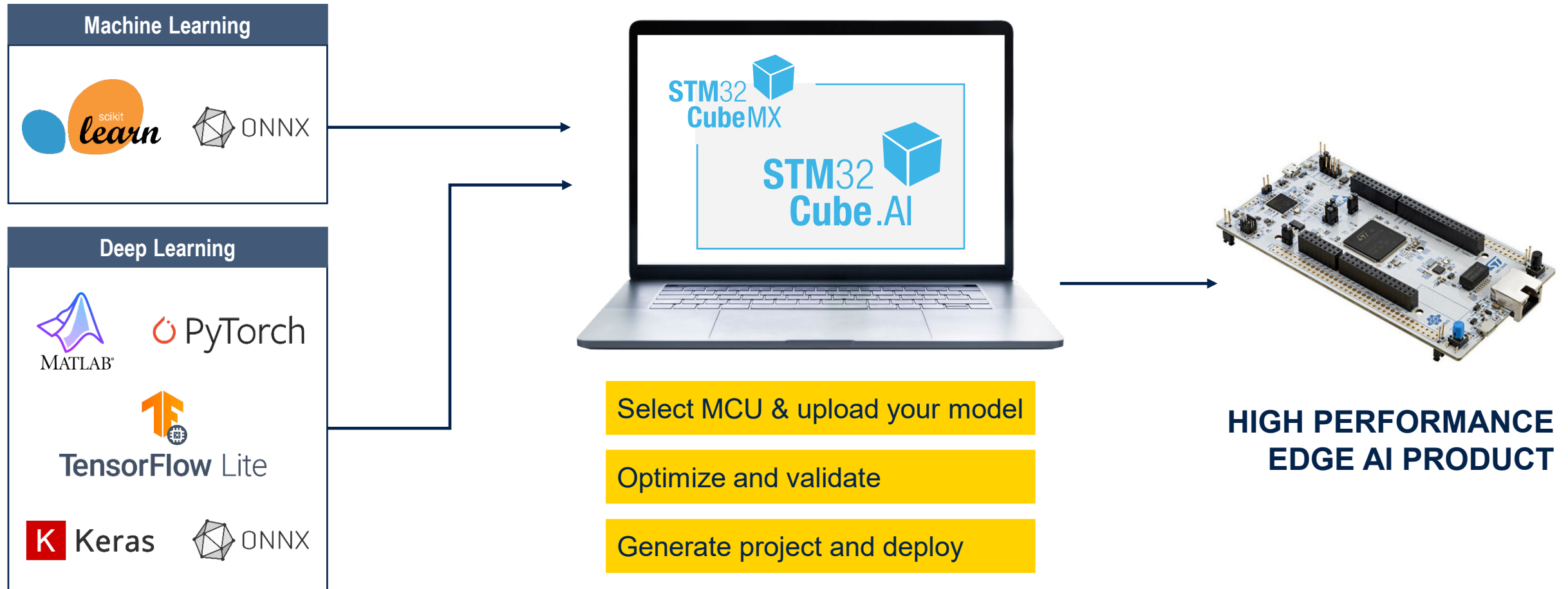No sharing in the cloud



**Privacy by design**
GDPR compliant

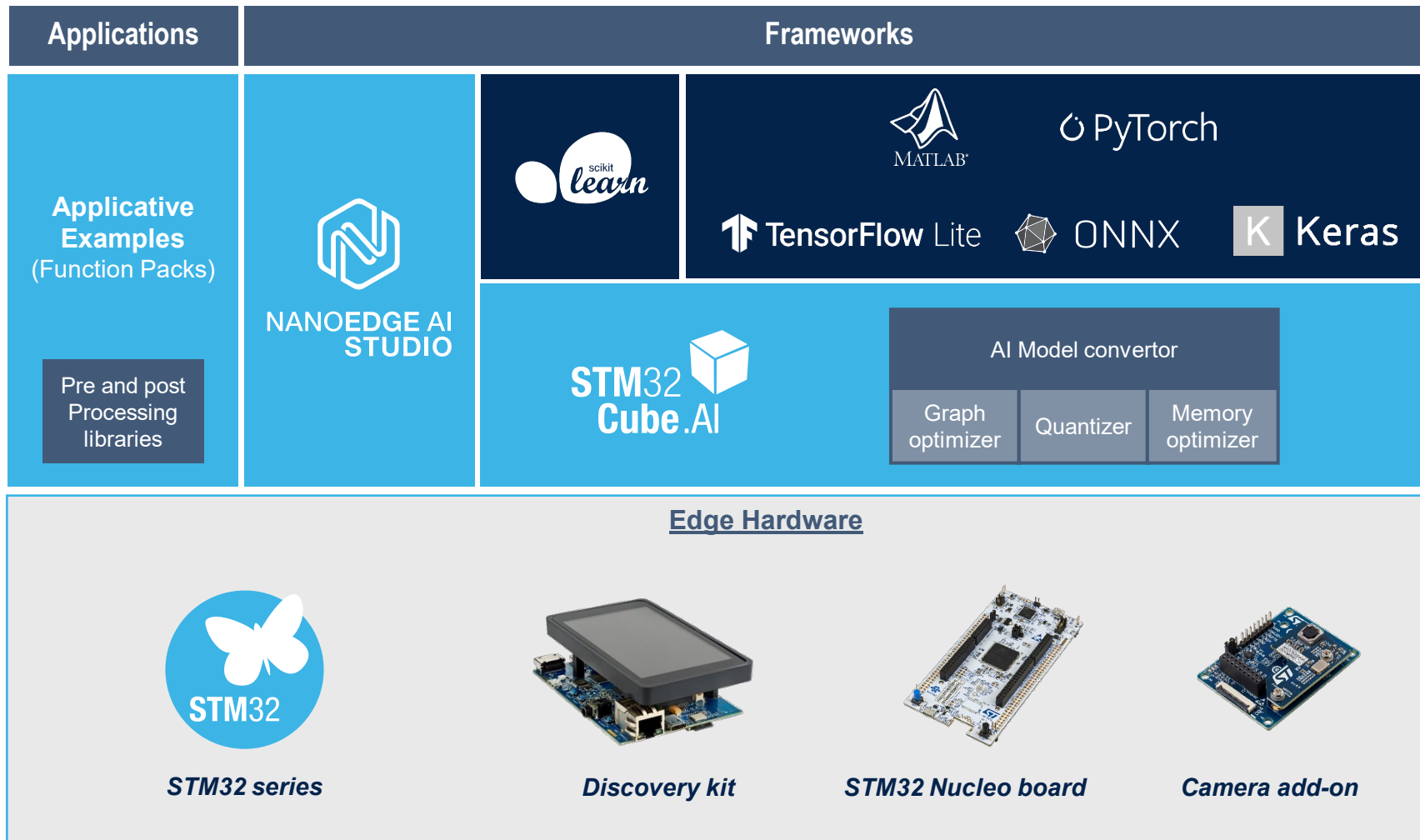**Sustainable on energy**
Low-power consumption

**Better user experience**

# STM32 Cube.AI

# A tool to seamlessly integrate AI in your projects

**Machine Learning**

scikit learn | ONNX

**Deep Learning**

MATLAB | PyTorch
TensorFlow Lite
Keras | ONNX

STM32 CubeMX
STM32 Cube.AI

Select MCU & upload your model

Optimize and validate

Generate project and deploy

**HIGH PERFORMANCE EDGE AI PRODUCT**

# STM32 comprehensive AI ecosystem

| Applications | Frameworks | | |
|---|---|---|---|
| **Applicative Examples** (Function Packs) | **NANOEDGE AI STUDIO** | scikit learn | MATLAB    PyTorch    TensorFlow Lite    ONNX    K Keras |
| Pre and post Processing libraries | | **STM32 Cube.AI** | AI Model convertor — Graph optimizer / Quantizer / Memory optimizer |

**Edge Hardware**



**STM32 series**          **Discovery kit**          **STM32 Nucleo board**          **Camera add-on**

# The 3 pillars of STM32Cube.AI

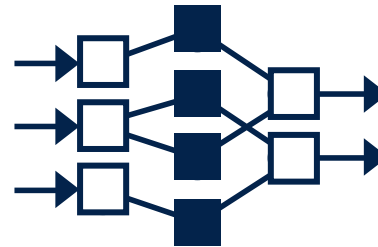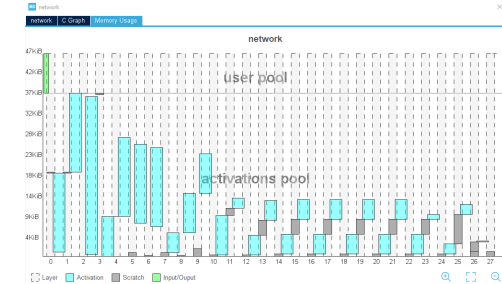| Graph optimizer | Quantized model support | Memory optimizer |
|---|---|---|
| Automatically improve performance through graph simplifications & optimizations that benefit STM32 target HW architectures | Import your quantized ANN to be compatible with STM32 embedded architectures while keeping their performance | Optimize memory allocation to get the best performance while respecting the constraints of your embedded design |



- Auto graph rewrite
- Node/operator fusion
- Layout optimization
- Constant-folding…
- Operator-level info to fine-tune memory footprint and computation

- From FP32 to Int8
- Minimum loss of accuracy
- Code validation on target
  o Latency
  o Accuracy
  o Memory usage

- Memory allocation
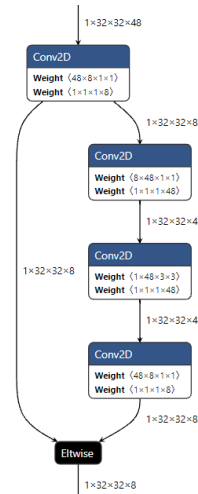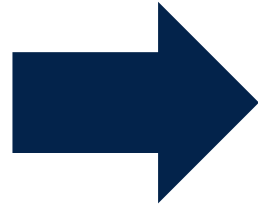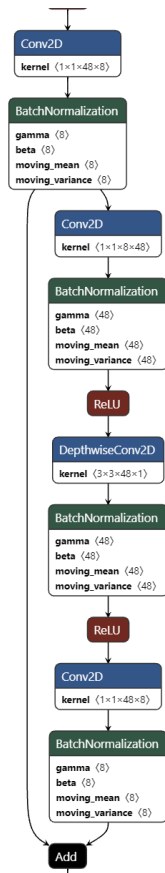- Internal/external memory repartition
- Model-only update option

STM32Cube.AI is **free of charge**, available both in graphical interface and in command line.
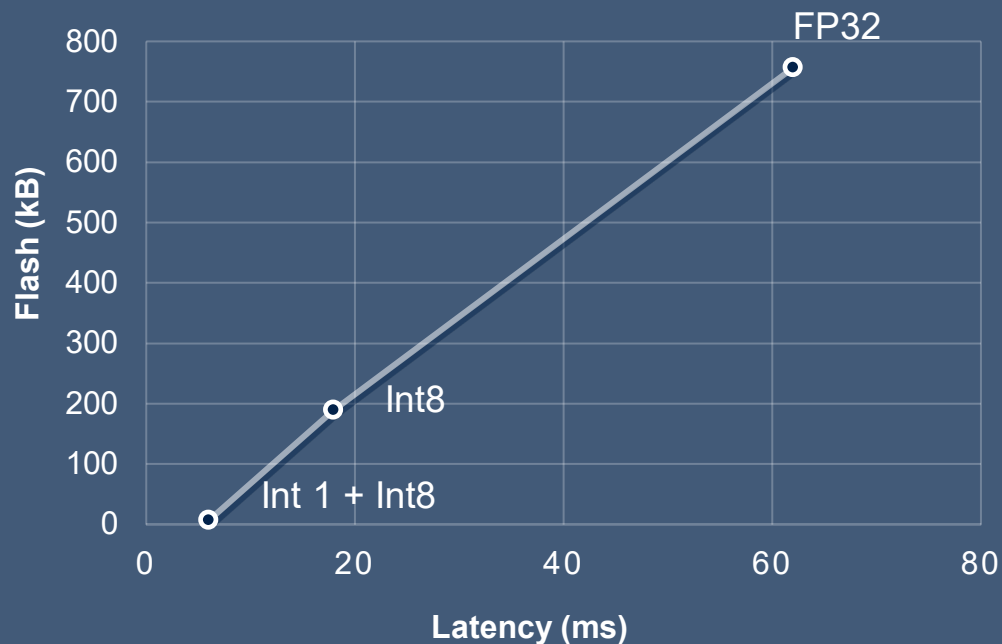
## Squeeze your graph to fit into an MCU!



**Fully automated process in the STM32Cube.AI workflow**

- Your original graph is optimized at the very early stage for optimal integration into STM32 MCU/MPU

- Loss-less conversion

# Quantized model support

**Simply use quantized networks to reduce memory footprint and inference time**

**LATENCY & MEMORY COMPARISON FOR QUANTIZED MODELS**

[Chart: Flash (kB) vs Latency (ms)]
- FP32 (~62 ms, ~760 kB)
- Int8 (~18 ms, ~190 kB)
- Int 1 + Int8 (~5 ms, ~10 kB)

Flash (kB): 0, 100, 200, 300, 400, 500, 600, 700, 800
Latency (ms): 0, 20, 40, 60, 80

STM32Cube.AI support quantized Neural Network models with **all parameter formats**:
- FP32
- Int8
- Mixed binary Int1 to Int8 (Qkeras*, Larq.dev*)

*Please contact edge.ai@st.com to request the relevant version of STM32Cube.AI*
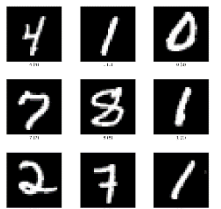
**HW Target**: NUCLEO-STM32H743ZI2
**Model**: Low complexity handwritten digit reading
**Freq**: 480 MHz
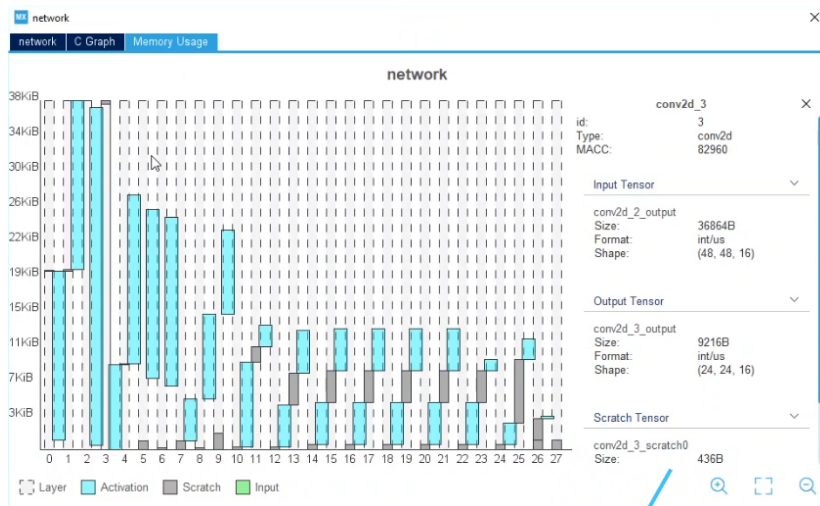**Accuracy**: >97% for all quantized models

**Tested database**: MNIST dataset

MNIST dataset

## Optimize performance easily with the memory allocation tool



**Model memory allocation**

- Set your external memory
- Map in non-contiguous internal flash section
- Partition internal vs external flash memories

**Model RAM consumption per layer**

- Easily identify most critical layers

**Re-use model input buffer to store activation data***

- Minimize RAM requirements

**Relocatable network**

- A separate binary is generated for the library and the network to enable standalone model upgrade

*\* Requires input and activation buffers in same memory*

# Making Edge AI accessible to all STM32 portfolio

**STM32Cube.AI is compatible with all STM32 series**

| | MPU | | | | | STM32**MP1** |
|---|---|---|---|---|---|---|
| | | | | | | 4158 CoreMark |
| | | | | | | Up to 800 MHz Cortex –A7 |
| | | | | | | 209 MHz Cortex –M4 |

| | High Perf MCUs | STM32**F3** | STM32**G4** | STM32**F2** | STM32**F4** | STM32**F7** | STM32**H7** |
|---|---|---|---|---|---|---|---|
| | | 245 CoreMark | 569 CoreMark | Up to 398 CoreMark | Up to 608 CoreMark | 1082 CoreMark | Up to 3224 CoreMark |
| | | 72 MHz Cortex-M4 | 170 MHz Cortex-M4 | 120 MHz Cortex-M3 | 180 MHz Cortex-M4 | 216 MHz Cortex-M7 | Up to 550 MHz Cortex -M7 |
| | | | | | | | 240 MHz Cortex -M4 |

Optimized for mixed-signal Applications

| | Mainstream MCUs | STM32**F0** | STM32**G0** | STM32**F1** |
|---|---|---|---|---|
| | | 106 CoreMark | 142 CoreMark | 177 CoreMark |
| | | 48 MHz Cortex-M0 | 64 MHz Cortex-M0+ | 72 MHz Cortex-M3 |

| | Ultra-low Power MCUs | STM32**L0** | STM32**L1** | STM32**L4** | STM32**L4+** | STM32**L5** | STM32**U5** |
|---|---|---|---|---|---|---|---|
| | | 75 CoreMark | 93 CoreMark | 273 CoreMark | 409 CoreMark | 443 CoreMark | 651 CoreMark |
| | | 32 MHz Cortex-M0+ | 32 MHz Cortex-M3 | 80 MHz Cortex-M4 | 120 MHz Cortex-M4 | 110 MHz Cortex-M33 | 160 MHz Cortex-M33 |

| | Wireless MCUs | | STM32**WL** | STM32**WB** | |
|---|---|---|---|---|---|
| | | | 162 CoreMark | 216 CoreMark | |
| | | | 48 MHz Cortex-M4 | 64 MHz Cortex-M4 | |
| | | | 48 MHz Cortex-M0+ | 32 MHz Cortex-M0+ | |

Latest product generation

STM32Cube.AI
The STM32CubeMX expansion pack for ML

Power Consumption Calculator

MCU Selector

Code Generation

Pinout Configuration

Middleware Parameters

Clock Tree Initialization

Peripherals Configuration

16

# Integrate your ML models more easily with our application-oriented code examples

## Time series-based monitoring



### FP-AI-MONITOR1

- Predictive maintenance and much more sensor-monitoring apps
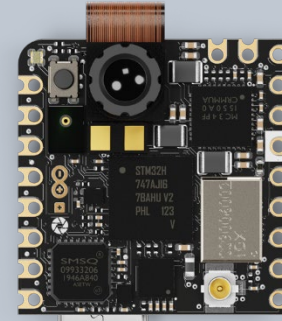- Runs Libraries from NanoEdge™ AI Studio

## Audio and Sensing



### FP-AI-SENSING1

- Human Activity Recognition
- Acoustic Scene Classification
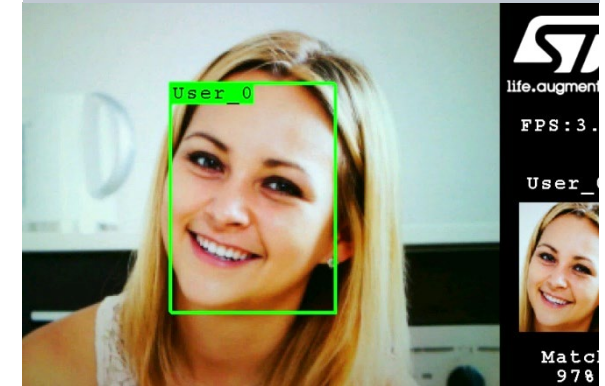- Data logging, labeling and result on BLE applications

## Computer Vision



### FP-AI-VISION1

- Food recognition (CNN)
- Person presence detection (CNN)
- People counting (Object detection NN)
- Image processing Library
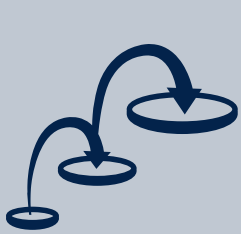
## Face recognition



### FP-AI-FACEREC1

- Face detection and recognition
- Fully functional without cloud connection

# We provide everything to kick off your project

## Design documentation



**Getting started**
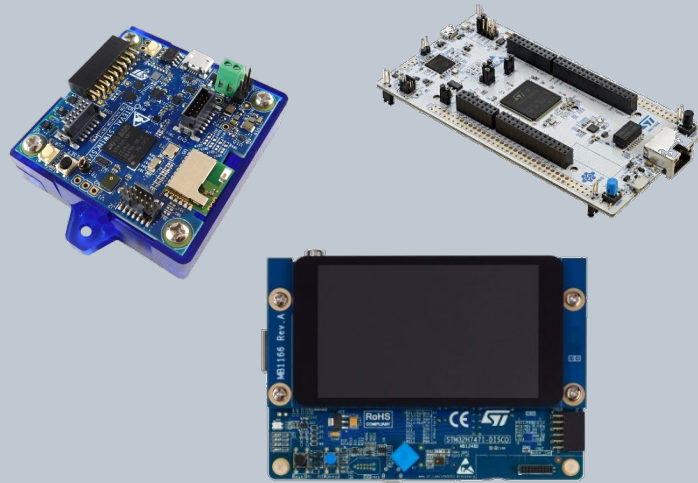
Be guided step-by-step to learn STM32 ecosystem

**Development zone**

Get started on application development and project sharing

- **Wiki by ST** is a great forum to learn and start developing AI on STM32!
- Videos of application examples
- Massive Open Online Course (MOOC)

## Hardware and software tools



- Evaluation platforms for STM32 MCU/MPU
- Extra sensor boards
- Full software suite

## Support & Updates



- **ST Community**: STM32 ML & AI group
- Distributor certified FAE
- Support center
- Newsletter

# What's new in STM32Cube.AI v7.1.0?

**v7.1.0**

**Bringing STM32Cube.AI to all STM32 and improving overall performances**

### # Now supporting entry level STM32

- Introducing the support for STM32 arm Cortex-M0 and arm Cortex-M0+
- STM32Cube.AI can now generate optimized code for STM32C0, STM32F0, STM32L0 and STM32G0 series

### # Improved user experience and performance tuning

- Added advanced support for splitting the activation buffer over several memory segments (multi-heap support) allowing full manipulation of the different on-board memories of the STM32H7 for example.

### # Up-to-date and improved code generation

- Support TensorFlow Lite micro v2.7 runtime and ONNX 1.9
- Support of more Keras, TensorFlow Lite and ONNX layers (refer to documentation for exhaustive list)
- Extend support of scikit-learn algorithms with new ONNX-ML operators

# Don't go alone

## We have created a network of companies to support you

Trust our **authorized partners** to ensure the success of your project. Learn more at st.com/stm32ai

Wish to discuss a co-development partnership for ML/AI projects? Contact us at edge.ai@st.com

# Releasing your creativity



f  /STM32

🐦 @ST_World

🔗 community.st.com

🌐 www.st.com/STM32ai

Wiki wiki.st.com/stm32

🐙 github.com/STMicroelectronics

▶ Videos

📰 STM32Cube.AI blog articles

STM32

life.augmented

# Our technology starts with You

🌐 Find out more at www.st.com/stm32ai

life.augmented