



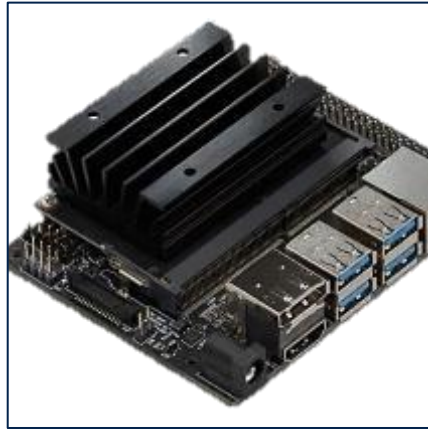
life.augmented

# STM32Cube.AI Neural Networks on STM32

Blaine Moon



# AI on the news...



**Can you easily access it?**



**What if we talk about bringing AI to your 2020 projects?**

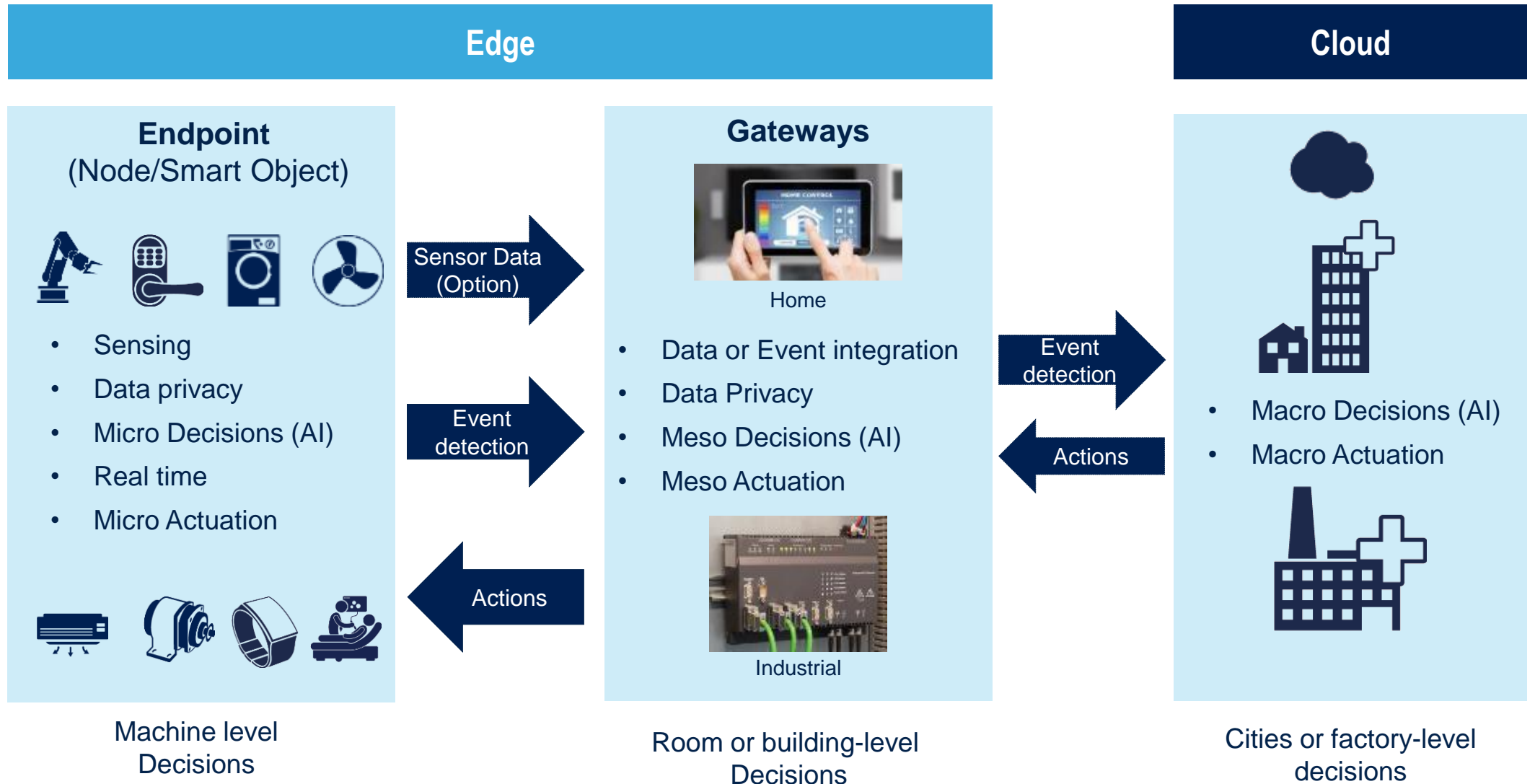
# ST and AI on the Edge



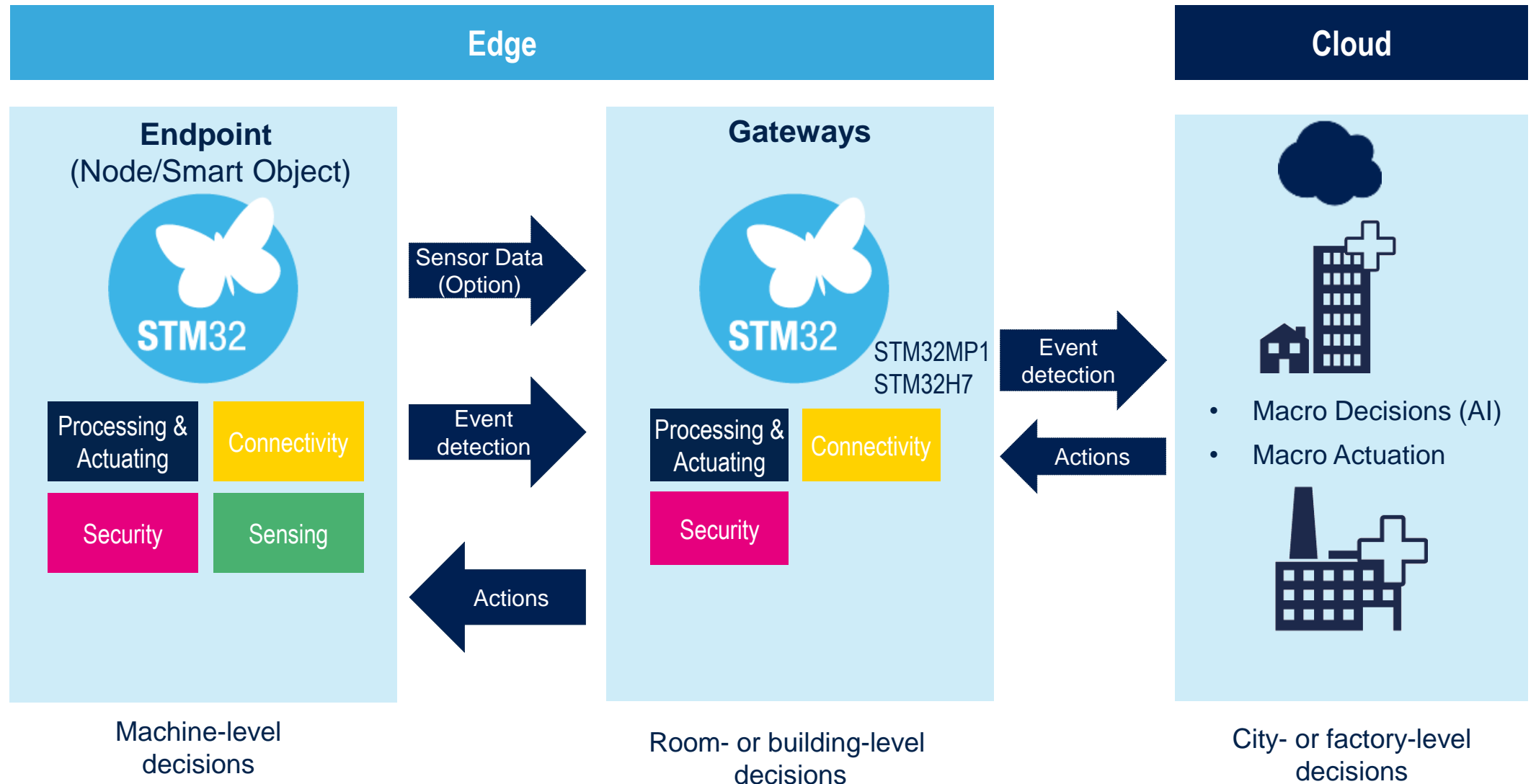
- **Microcontrollers** are the brain of these smart objects
- AI is a disruptive technology where ST has been investing for many years ...
- Local processing to fix limitations :
  - Latency and cost communications
  - Autonomy (Battery-operated devices)
  - Limited networks bandwidth or connectivity loss
  - Data privacy



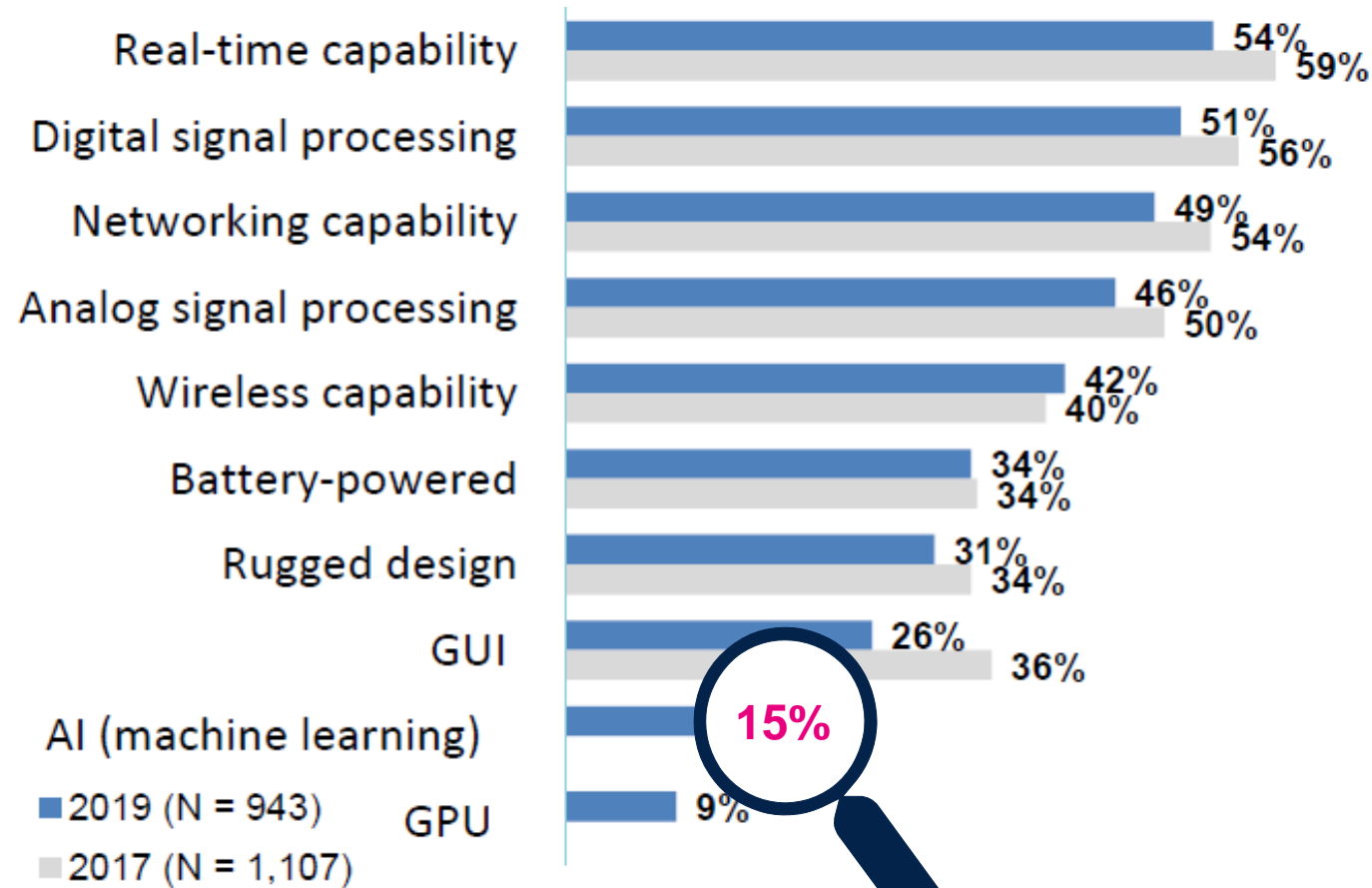
# Distributed AI from Edge to Cloud



# Distributed AI from Edge to Cloud

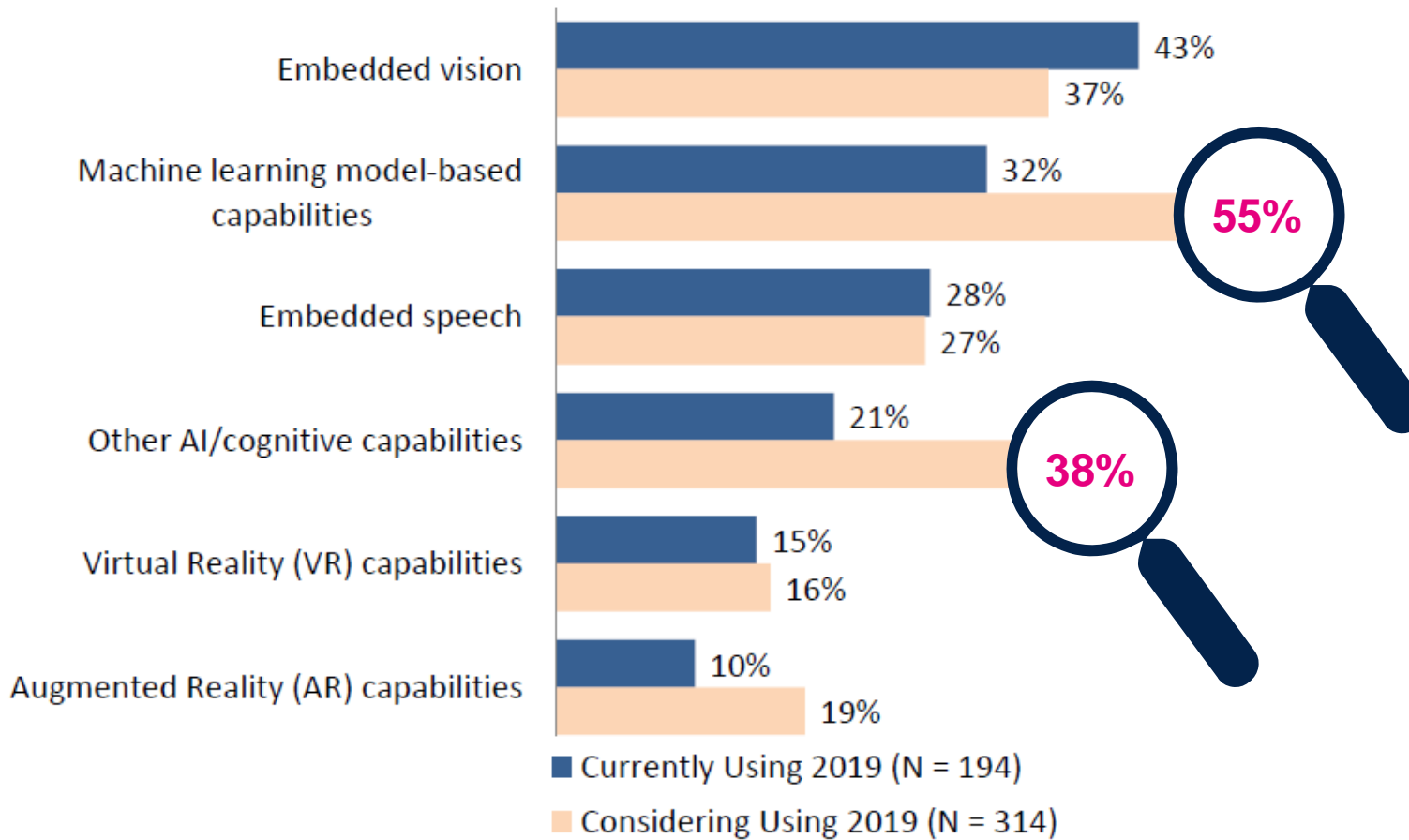


# AI moving to the Edge



- 15% of embedded projects already include AI in 2019
- 0% in 2017

# AI moving to the Edge



- MCU pervasion will grow
  - With Machine Learning
  - Other AI cognitive capabilities
- MPU adoption will continue
  - With Embedded vision & Speech

68%  
of EMEA users  
are considering  
using Machine  
Learninf

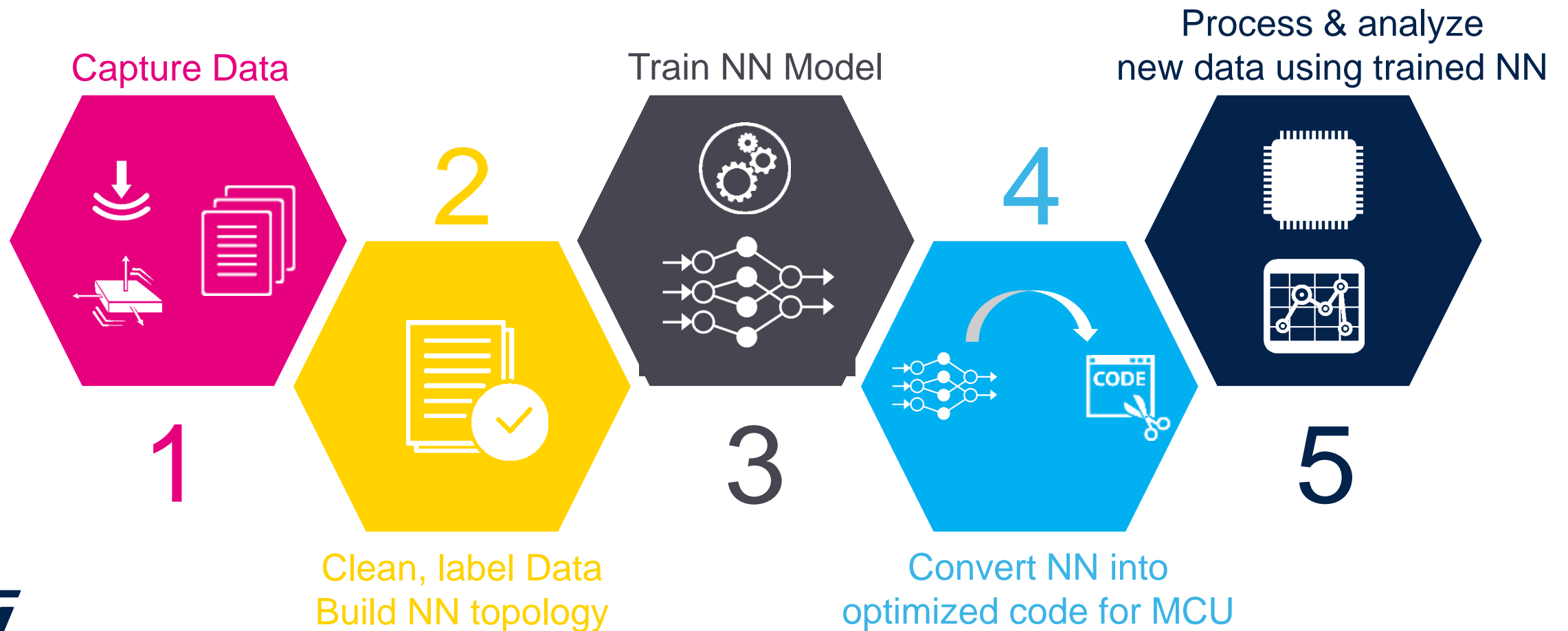
# Neural Networks on STM32

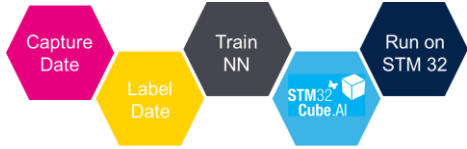
STM32    
Cube.AI



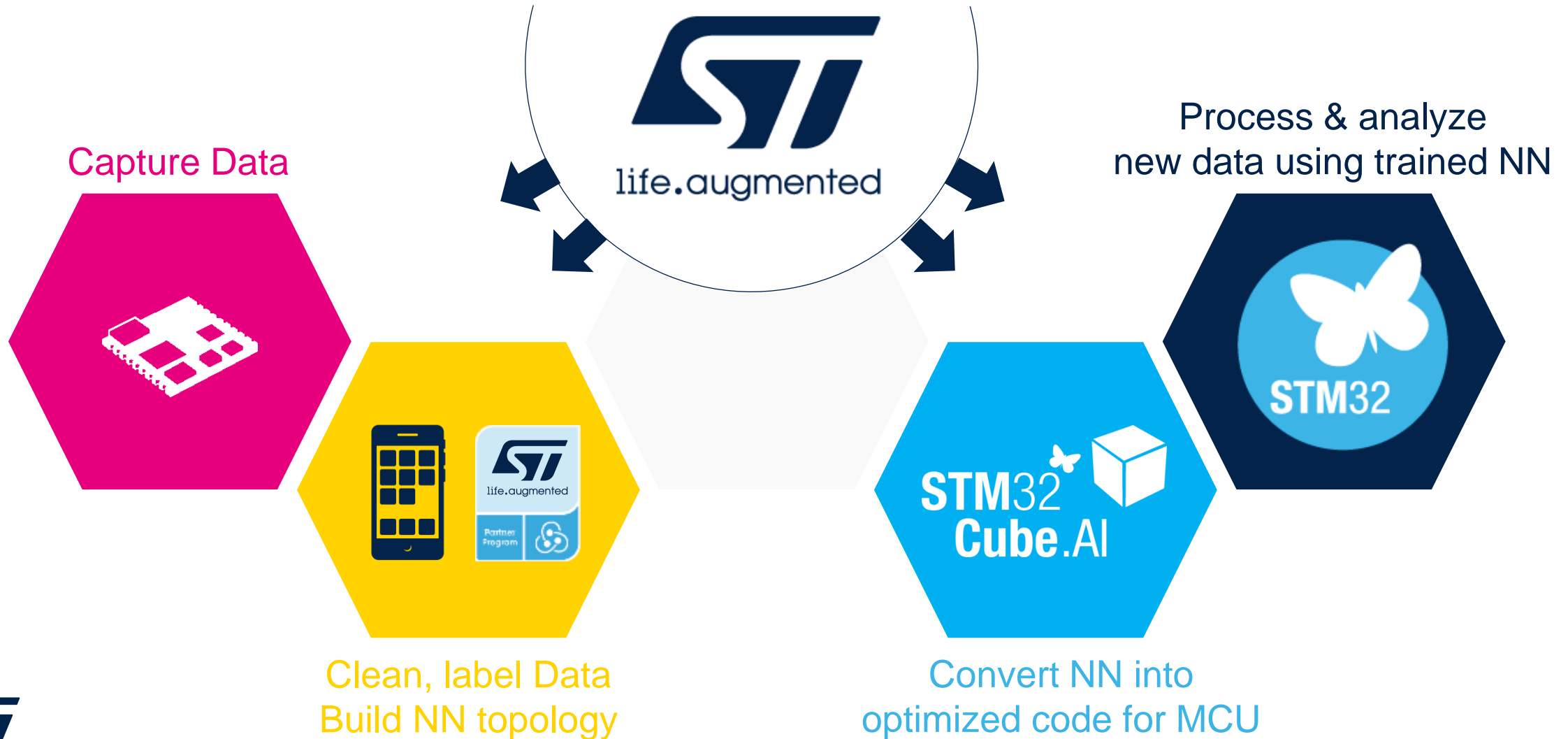


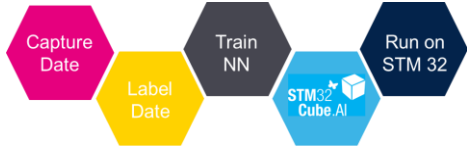
# The key steps behind Neural Networks





# ST toolbox for Neural Networks





# Collecting data & architecting an NN topology

## Services provided by Partners

## ST tools to support

### Capture Data



Clean, label Data  
Build NN topology



### ST BLE Sensor mobile phone application

Collect and label data from the SensorTile.



ST BLE  
Sensor



life.augmented

Partner  
Program



### Selected partners

Neural Networks engineering services support.  
Data scientists and Neural network architects.

# ST AI official partners

The screenshot shows the ST website header with the logo 'life.augmented' and navigation links: Products, Applications, Tools & Software, About ST, Sample & Buy, and Support & Community. Below the header is a search bar and a 'Contact Us' button. The main content area features 'STM32 solutions for Artificial Neural Networks' and a navigation bar with 'Overview', '5 Steps to AI', 'Resources', and 'Featured Content'. Under 'Resources', there are icons for 'Function Pack', 'Videos', and 'Dev Kit'. The 'Find a Partner' link is circled in pink, with a pink arrow pointing from it to the partner search interface on the right.

## Find a partner

Products and Services

- Cloud
- Embedded Software
- Software Dev Tools
- Components & Modules
- Engineering Services
- Training
- Hardware Dev Tools





STM32CubeAI

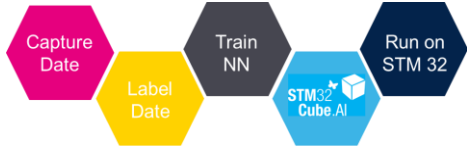
Country of operation

RESET Partner Search

Total results:

20 partners per page

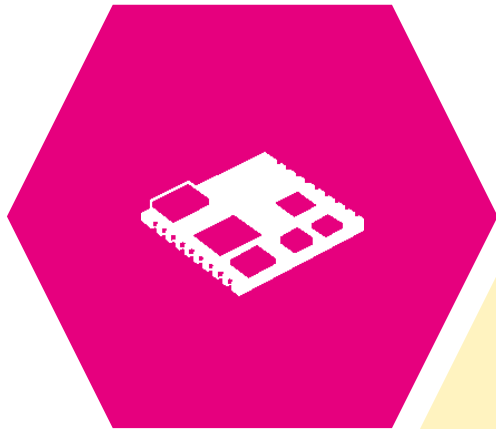
-  **Bluewind**  
Bluewind, an independent engineering company, provides innovative product design solutions in the fields of Electronics, Energy Efficiency, and Connected Devices. The R&D task force consists of 20+ experienced engineers. [Show more >](#)
-  **Cartesiam**  
Cartesiam is an expert in Artificial Intelligence at the Edge. We are an ISV expert in mathematics, AI and signal processing on microcontrollers. Cartesiam, invented NanoEdge™, a revolutionary technology enabling Machin... [Show more >](#)
-  **Imaginob**  
Imaginob is a global leader in artificial intelligence products for STM32. Based in Stockholm, Sweden, the company has been serving customers within the automotive, manufacturing, healthcare and lifestyle industries sin... [Show more >](#)
-  **Inventhys**  
Inventhys is a fast-growing, end-to-end IoT development services company. Our engineers design hardware, embedded



# Example form-factor hardware to capture and process data

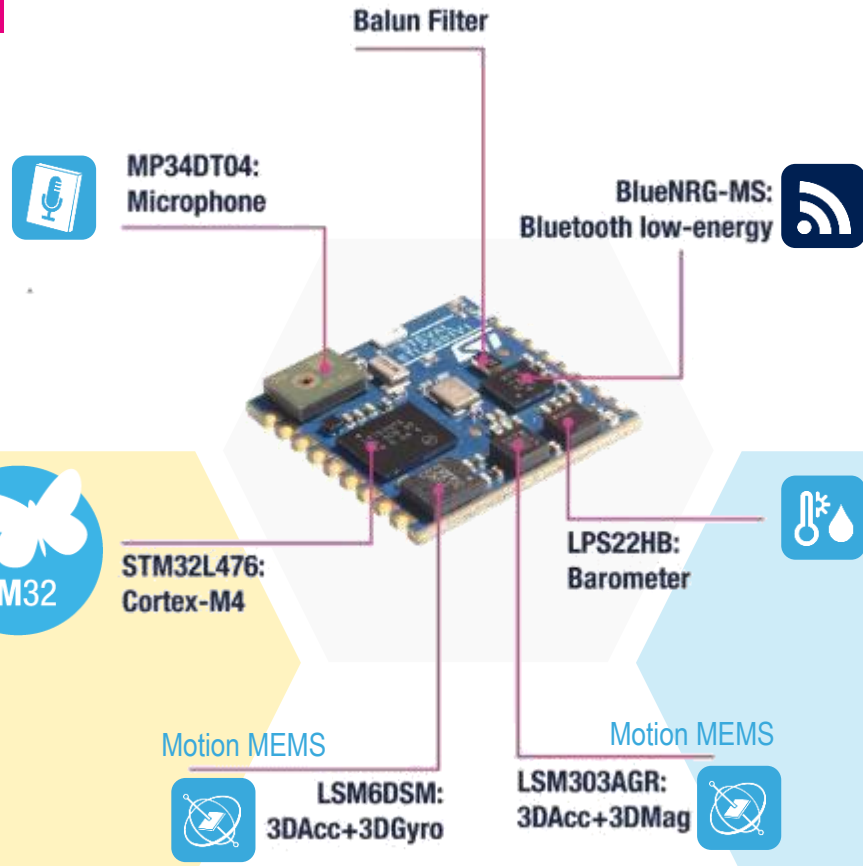
## SensorTile

Capture Data

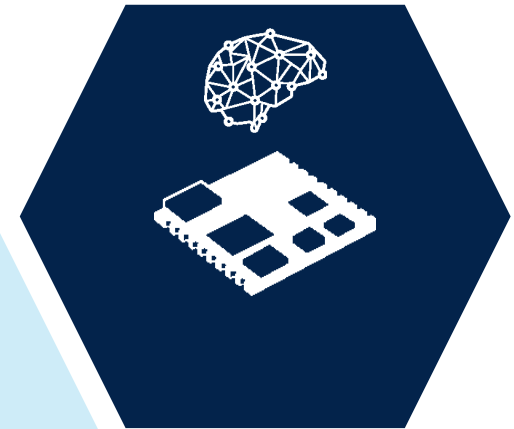


STM32

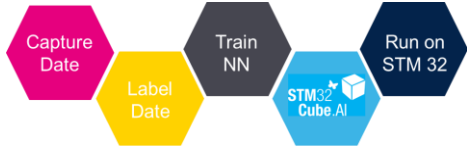
STM32L476: Cortex-M4



Process & analyze new data using trained NN



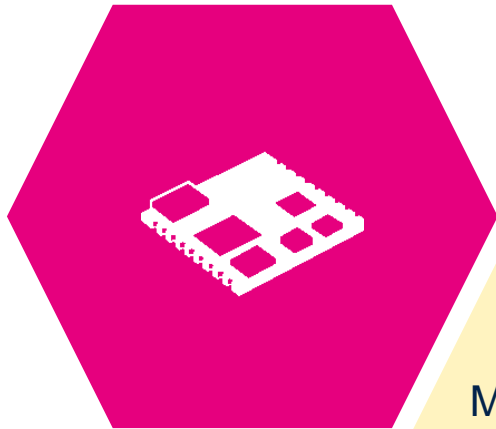




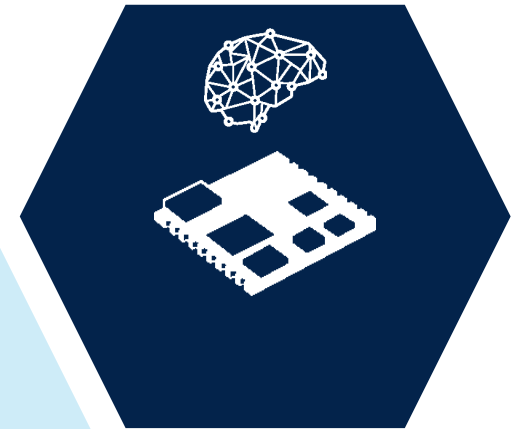
# Fast go-to-market module to capture data with more accuracy

## SensorTile.Box

Capture Data



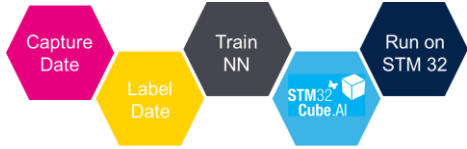
Process & analyze new data using trained NN



Microsoft IoT Services ready  
Microsoft Azure

More advanced, high accuracy and low power sensors

- First Inertial module with Machine Learning capabilities.
- Motion (accelerometer and gyroscope, magnetometer) and slow motion (inclinometer)
- Altitude (pressure), environment (pressure, temperature, humidity, compass) and sound (sound and ultrasound analog microphone)
- Microsoft IoT services ready to make available on a web dashboard the result of the embedded processing

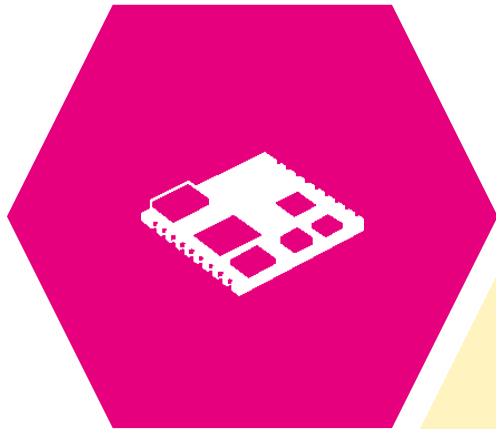


# Form factor hardware AI IoT node for more connectivity

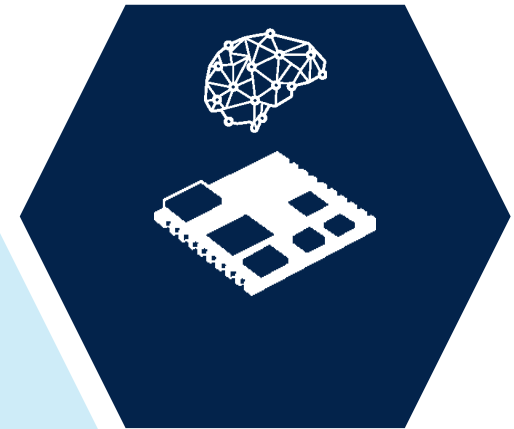
IoTNode

- +
-  Sub-1GHz
  -  Dynamic NFC Tag
  -  Wi-Fi

Capture Data

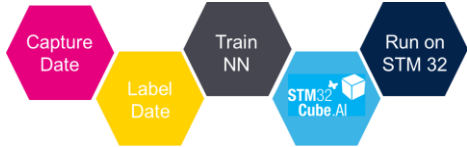


Process & analyze  
new data using trained NN



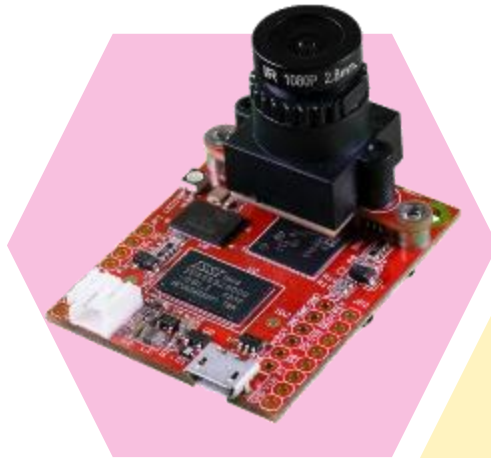
More debug capabilities

- Integrated ST-Link/V2.1
- PMOD extension connector
- Arduino Uno extension connectors



# OpenMV integration

## Fast machine vision prototyping



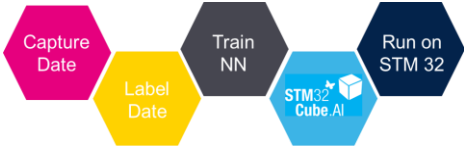
OpenMV CAM  
Running MicroPython on STM32

Configure Machine Vision in real-time over USB in Python



Run and validate optimized Neural Network



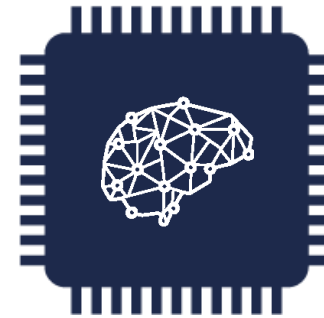
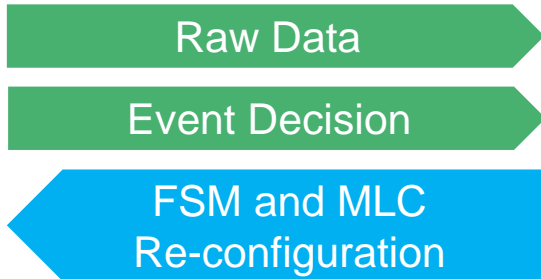


# Distributed AI: Sensor + STM32

## Optimize performance & power consumption

### Smart Sensor with Machine Learning Core

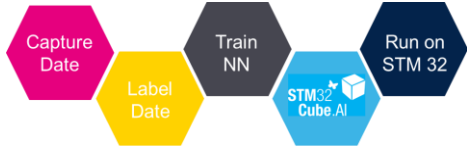
### Smart STM32 second level of AI processing



**Deep Learning  
Neural Networks  
Machine Learning**

- Best ultra-low-power sensing at high performance
  - 550µA (gyroscope and accelerometer)  
-25% power compared to competition
  - 20~40µA (Accelerometer only for HAR)
- Efficient Finite State Machines: 2µA
- Configurable Machine Learning Core: 4~8µA

- More advanced and complex Neural Networks
- Input can be data from multiple sensor data and/or sensor Machine Learning decisions
- Possible to run multiple Neural Networks
- Actuation & communication



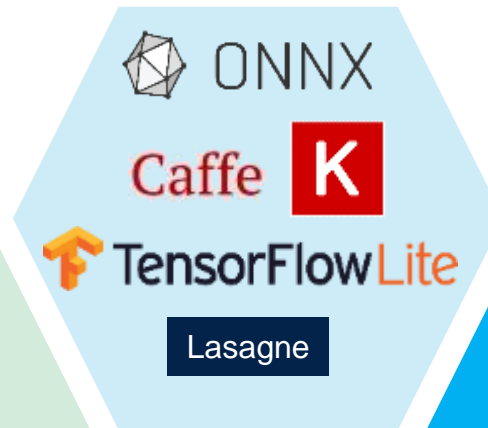
# STM32CubeMX extension AI conversion tool

Input your framework-dependent, pre-trained Neural Network into the **STM32Cube.AI** conversion tool

Automatic and fast generation of an STM32-optimized library

**STM32Cube.AI** offers interoperability with state-of-the-art Deep Learning design frameworks

Train NN Model

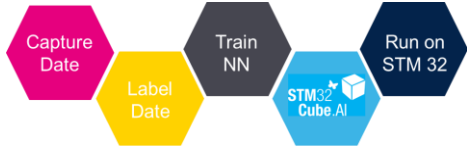


Process & analyze new data using trained NN



Convert NN into optimized code for MCU  
Validate code directly on target





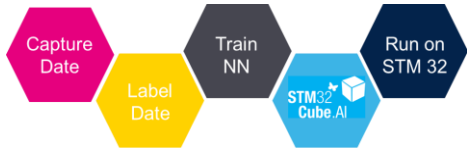
# Import Neural Networks

## Open Neural Network (ONNX) exchange format



- ▶ ONNX open format enables models to be trained in one framework and transferred to another for inference
- ▶ Common import/export format of many frameworks
- ▶ STM32 Cube.AI hardware optimization is available for any tools exporting ONNX models





# STM32 solutions for AI

## More than just the STM32Cube.AI

An extensive toolbox to support easy creation of your AI application



AI extension for STM32CubeMX to map pre-trained Neural Networks



Software examples for Quick prototyping  
Audio, Motion and Vision Function packs  
On **ST development Hardware**



STM32 **Community** with dedicated Neural Networks topic



Trainings, hands-on exercises, MOOCs and partners **videos**

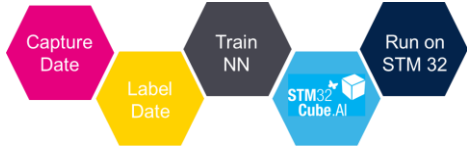


STM32 AI Partner Program with dedicated Partners providing **Machine or Deep Learning engineering services**

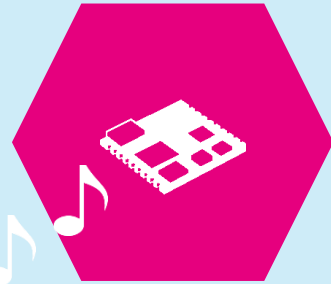
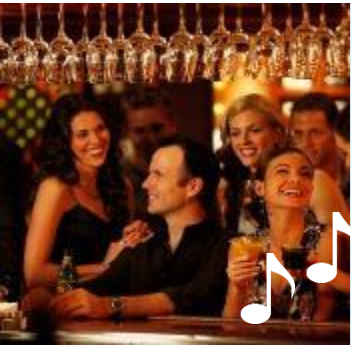
# Function Packs (Software examples)

STM32    
Cube.AI

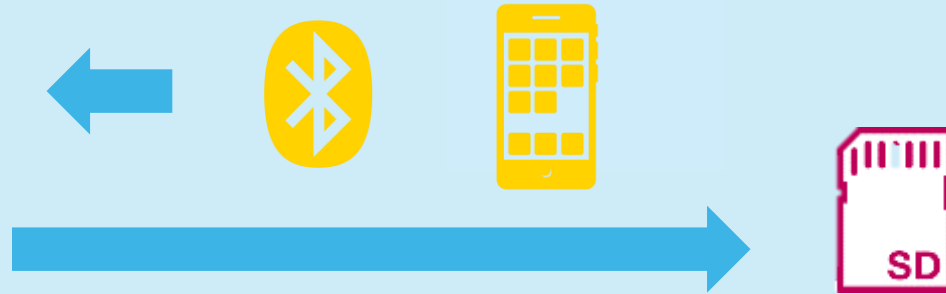




# Audio Scene Classification (ASC) Audio example in FP-AI-SENSING1 package



Audio Data capture



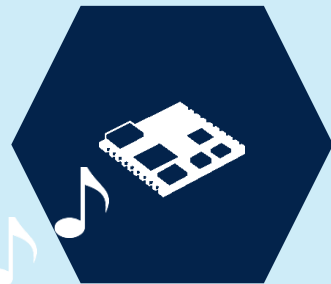
Labeling controlled by smartphone application

Data stored on the device SD card for future learning

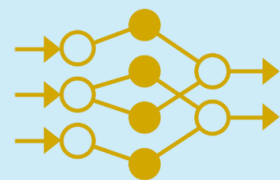


3 classes

Indoor, Outdoor, In vehicle labeling



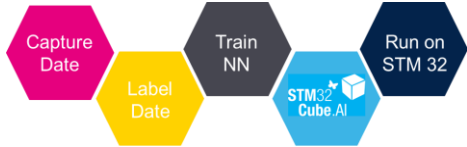
Embedded audio pre-processing



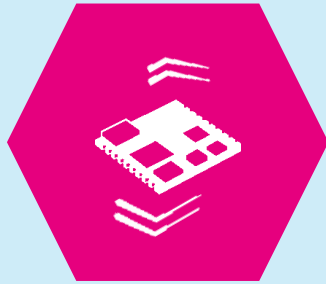
NN & example dataset provided



Inference result displayed on mobile app



# Human Activity Recognition (HAR) Motion example in FP-AI-SENSING1 package



**Motion Data Capture**



**Labeling** controlled  
by smartphone application

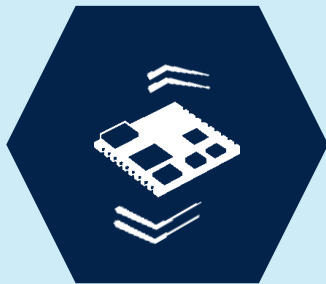


Data stored on the device  
SD card for future **learning**

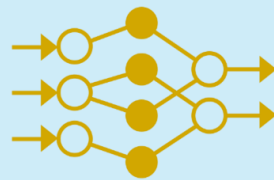


**5 classes example**

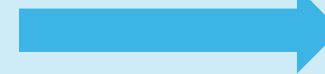
Stationary, walking, running,  
biking, driving **labeling**



**Embedded motion**  
pre-processing



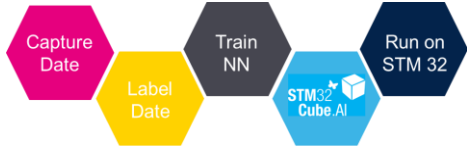
NN & example  
dataset provided



**Inference result**  
displayed on mobile app

**Inferences** running  
on the microcontroller





# Image Classification Vision example in FP-AI-VISION1 package

## Enjoy the **food classification demo**

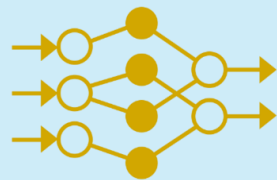
- Default demo based on 18 classes (224x224 RGB pictures)
- Several camera image output size possible

## Full **end-to-end optimized software example**

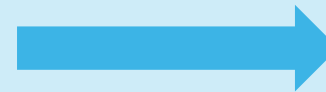
- from camera acquisition to image pre-processing before feeding the NN
- Multiple memory mapping possibilities to optimize and test impact on performances
- Retrain this NN with your own dataset
- Quantize your trained network to optimized inference time and memory usage



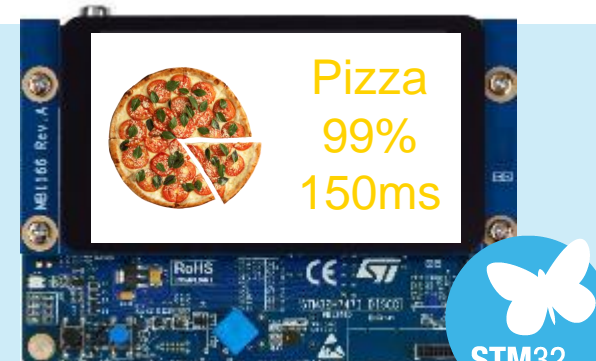
Embedded **image** pre-processing (SW) on the STM32H747



NN & example dataset provided



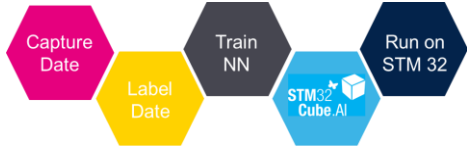
**Inferences** running on the microcontroller



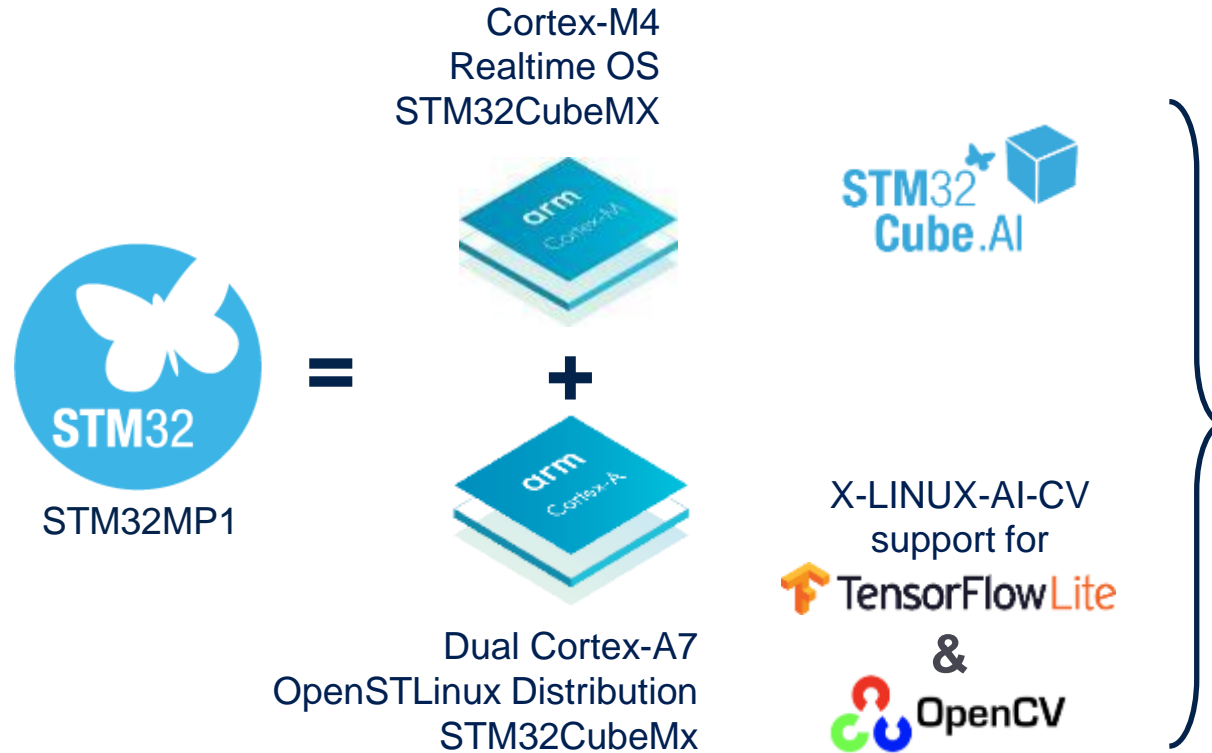
**Inference result** displayed on STM32H747 Discovery board LCD display

# AI solutions for STM32MP1

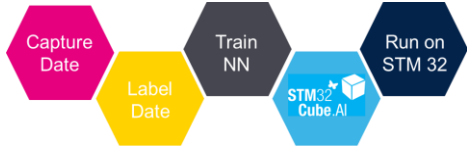




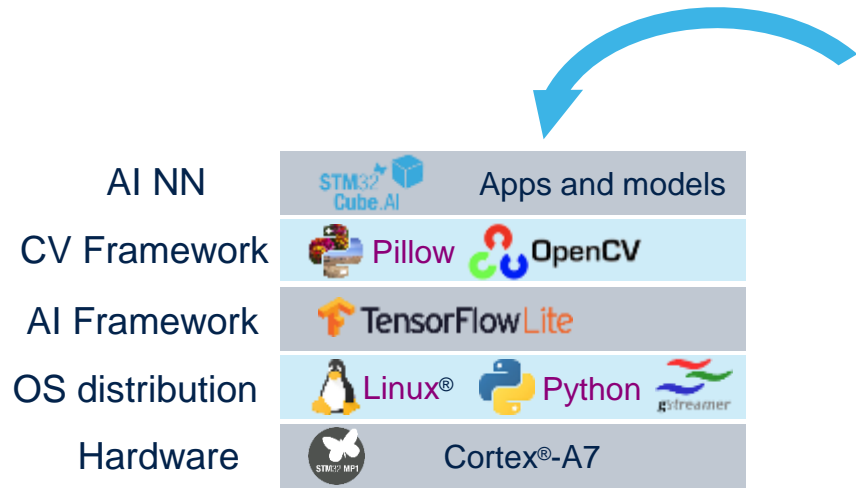
# STM32MP1 microprocessor Augmented intelligence



- STM32Cube.AI to convert pre-trained NNs for the Cortex-M4 core
- TensorFlow Lite STM32MP1 support up streamed for native NN inferences support on the dual Cortex-A side



# X-Linux-AI-CV package for STM32MP1 computer vision application



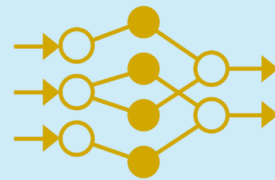
Application examples in C/C++ and Python

- Image classification: 1000 objects classified
- Multiple object detection: 90 classes

Includes code for camera acquisition and image pre-processing



USB camera or built-in camera module



AI, CV frameworks & application examples provided

**Inferences** running on the microprocessor in 80ms for image classification

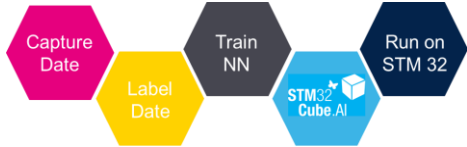


Displayed on STM32MP1-DK2, STM32MP1-EV1 and Avenger96 board

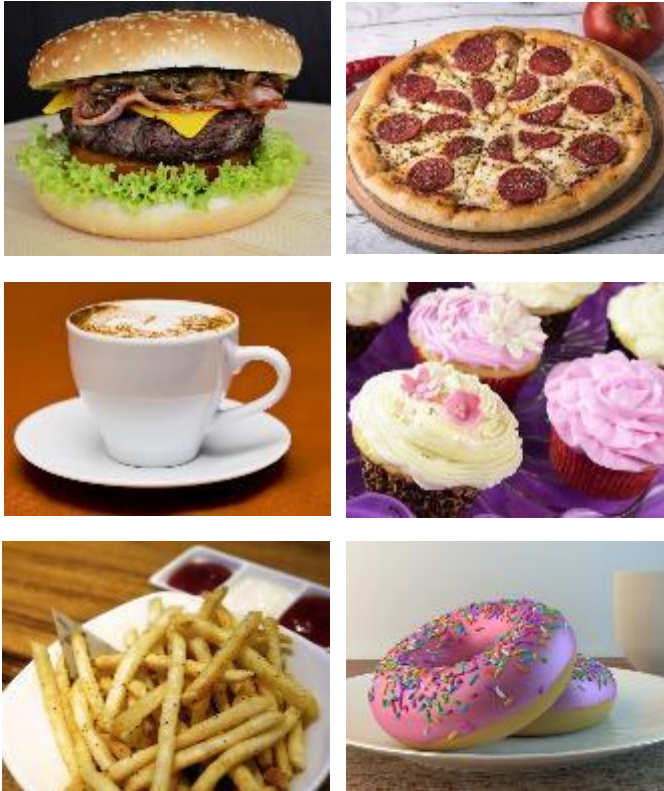
# Vision on MCU – A reality now diving into STM32Cube.AI







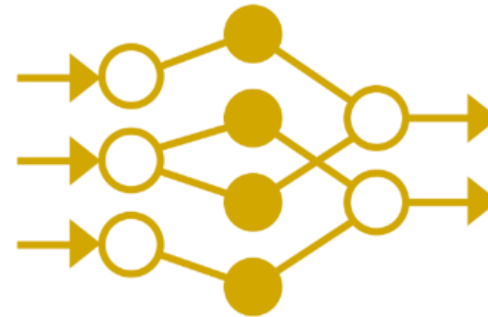
# Image Classification Vision example in FP-AI-VISION1 package



Live video stream



Embedded image pre-processing (SW) on STM32H747



**Inferences** running on the microcontroller

18 classes of food recognized



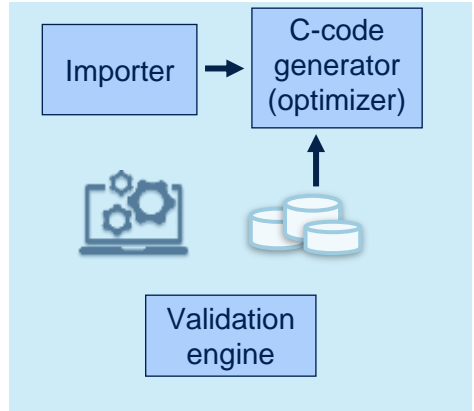
Recognition: Hamburger  
Confidence: 85%

**Inference result** displayed on STM32H747 Discovery board LCD



# The generation steps

- 1. Design the Neural Network
- ↓
- 2. Optimize the execution thanks to Cube.AI
- ↓
- 3. Validate the optimization with Cube.AI
- ↓
- 4. Integrate in a complete application thanks to function pack example



Function Pack



STM32H7

# STM32Cube.AI code generation

The screenshot displays the STM32Cube.AI software interface. The main window is titled "STM32CubeMX Untitled\*: STM32H743XIHx STM32H743I-EVAL". The interface is divided into several sections: "Pinout & Configuration", "Clock Configuration", and "Project Manager". The "Pinout & Configuration" section is active, showing the "Mode" and "Configuration" tabs. The "Configuration" tab is selected, and the "Main" network is highlighted. The "Model inputs" section shows a "Keras" model selected, with a "Saved model" dropdown menu. The "Model" field contains the path "d3/fd\_mobilenet\_food\_18\_mixed\_3\_original\_0.778.h5". The "Compression" is set to "None" and "Validation inputs" is set to "Random numbers".

On the right side, a "foodrecognition" window is open, showing a neural network diagram with the following layers:

- Input (ID: 0, Name: input\_2, Type: Input)
- Conv2D (ID: 1, Name: conv1, Type: Conv2D, Flash: 1792 B, MACC: 5619728)
- ScaleBias (ID: 2, Name: conv1\_bn, Type: ScaleBias)
- Nonlinearity (ID: 3, Name: conv1\_relu, Type: Nonlinearity)

At the bottom of the main window, the "Analysis" results are displayed:

|                  |   |
|------------------|---|
| macc             | : 24,985,408                                    |
| weights (rw)     | : 579,656 (566.07 KiB) (-0.86%)                 |
| activations (rw) | : 863,872 (843.62 KiB)                          |
| ram (total)      | : 1,466,056 (1.40 MiB) = 863,872 + 602,112 + 72 |

The "Analysis" status is "done" and the "Evaluation status" is "Acc RMSE MAE".

1 Enter the model:  
here float Keras  
h5 model

Memory footprint  
after analysis:  
> 1 MB of RAM

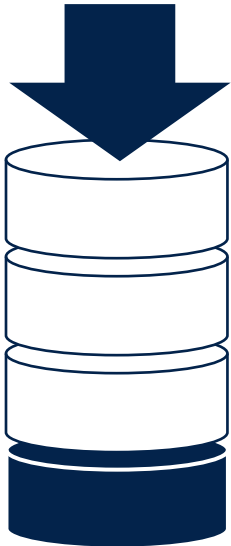
How to reduce?

2

3

# Quantization objective

Computer Vision Use Cases need a lot of memory  
→ Critical to reduce the memory for efficiency



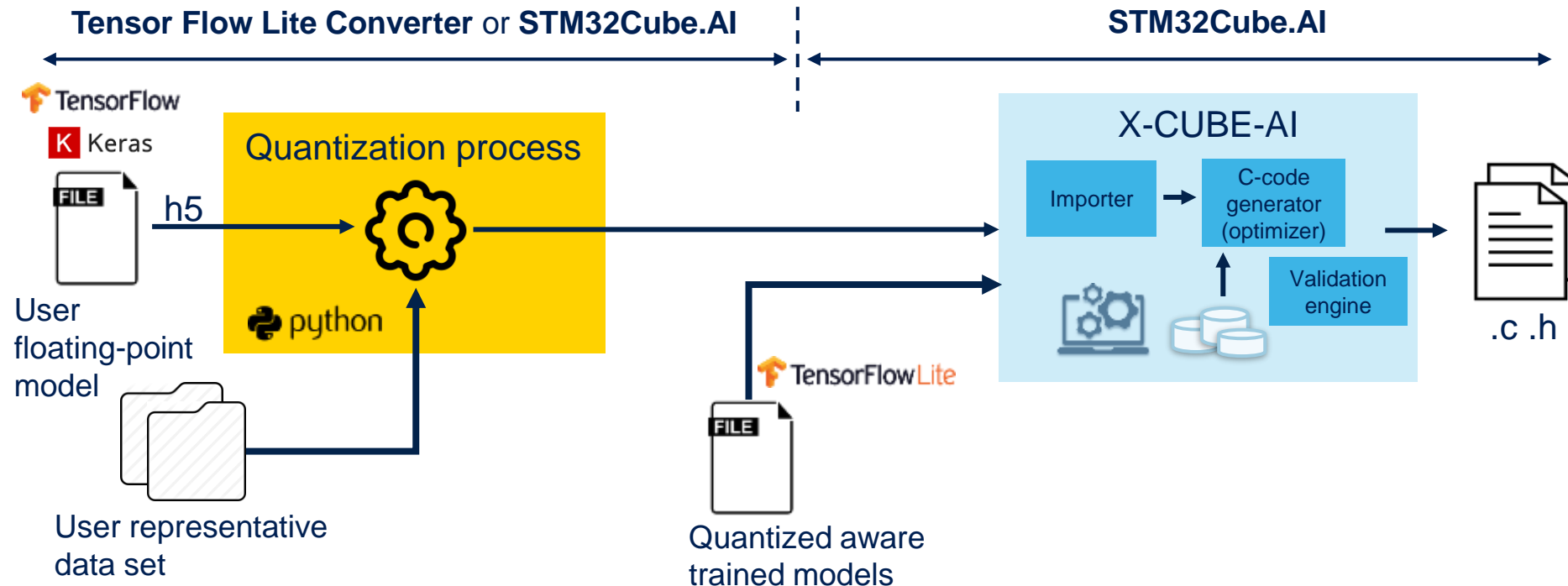
Convert floating-point (32 bits) model to quantized (8 bits) model

- **Reduce model size** (size of the memory to store the weights, in Flash). Up to x4
- **Reduce memory usage** (size of the activations buffer, in RAM). Up to x4
- **Improve latency**. Consequently **power consumption is improved**. Up to x3.
- With **minimal loss of accuracy**. Network size/complexity dependent

Reference: <STM32CubeMX\_localpath>\Packs\STMicroelectronics\X-CUBE-AI\5.0.0\Documentation

# STM32Cube.AI solution

- Post-training quantization algorithm
  - Based on an already-trained floating-point model
  - Parameters and activations are quantized
- Integrated code generation
  - Including adapted validation process
  - Support of TFL quantized aware trained models





# Results quantized vs float

Values measured on FP-AI-VISION1 1.1.0

**Using FD-MobileNet**  
(optimized 18-classes)

Input: 224x224x3  
MACC: 24 M  
Parameters: 145 K  
Nb of layers: 12  
Main kernel: DS-Conv,  
quantized 8 bits

| KPI                                | Float model           | Quantized model<br>Cube.AI 4.0.0 | Quantized model<br>Cube.AI 5.0.0 |
|------------------------------------|-----------------------|----------------------------------|----------------------------------|
| Flash                              | 566 kB                | <b>148 kB</b>                    | <b>148 kB</b>                    |
| RAM                                | 844 kB                | <b>393 kB</b>                    | <b>212 kB</b>                    |
| Accuracy                           | 77.8 %                | 77.1 %                           | 77.1 %                           |
| Inference time                     | 420 ms <sup>(1)</sup> | <b>155 ms</b>                    | <b>153 ms</b>                    |
| Frame Per<br>Second <sup>(2)</sup> | 2.2                   | <b>5.8</b>                       | <b>5.8</b>                       |

Inference times on the STM32H7 @400MHz, quantization done by Cube.AI (Qmn, greedy options)

<sup>(1)</sup> External memory: activation buffer & I/O buffers in external SDRAM

<sup>(2)</sup> FPS takes into account the capture, preprocessing and inference times.



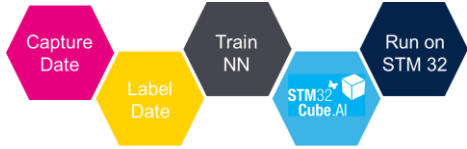
# Results through STM32 portfolio

## Performance measurement

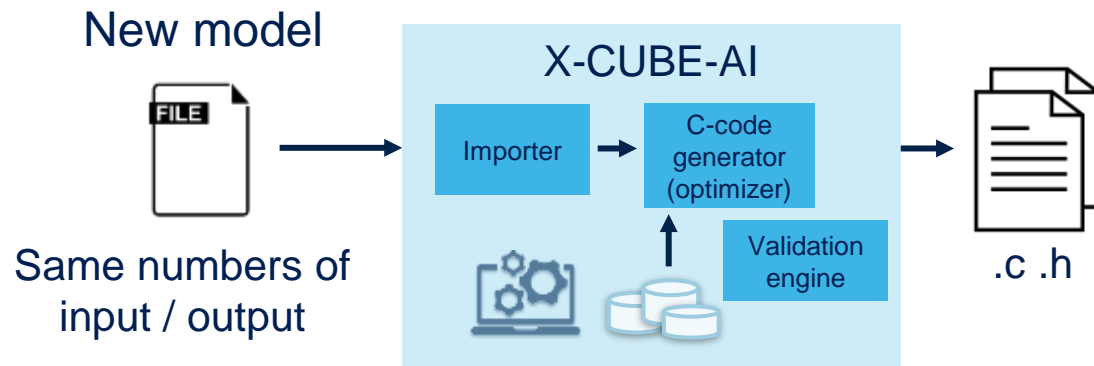
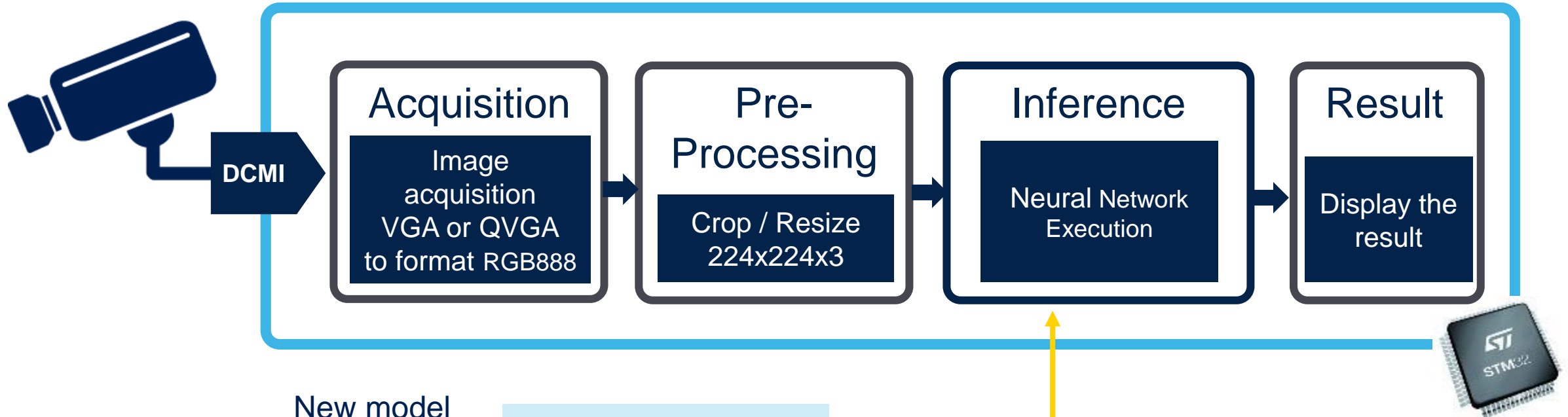
### Using FD-MobileNet (optimized 18-classes)

Input: 224x224x3  
 MACC: 24 M  
 Parameters: 145 K  
 Nb of layers: 12  
 Main kernel: DS-Conv,  
 quantized 8 bits

| KPI              | STM32L4R  | STM32H7                                   | STM32MP1                               |
|------------------|---|---|--|
| Flash            | 148 kB  | 148 kB                                    | 191kB                                  |
| RAM              | 212 kB  | 212 kB                                    | 1MB                                    |
| Inference Time   | 1.062 ms  | 153 ms                                    | 27.6 ms                                |
| Frame per Second | 0.9 fps   | 5.8 fps                                   | 36.2 fps                               |
| Processor        | L4R DK<br>M4@120 MHz<br>int RAM, int Flash,   | H747<br>M7@400 MHz<br>int RAM, int Flash, | MP1 DK2<br>2xA7@650 MHz<br>TFL 2.0 C++ |
| Cube.AI          | Cube.AI 5.0.0, 8 bits Ua/Ua quantized model<br>*input buffer allocated in activation buffer |   |  |

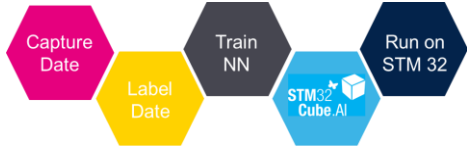


# FP-AI-VISION1 Model update

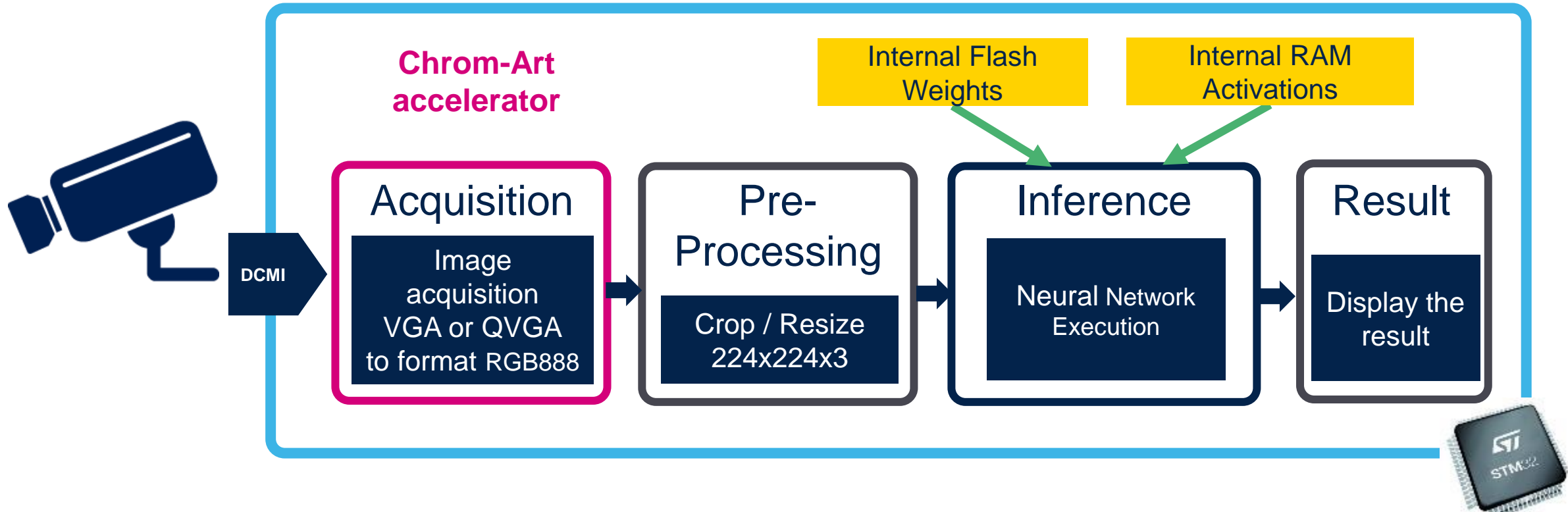


## Update only 4 files:

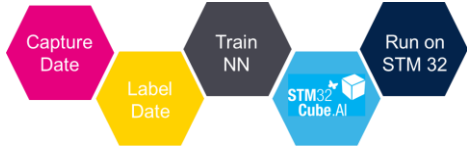
- network.h      public API
- network.c      network topology
- network\_data.h      parameters
- network\_data.c      weights



# FP-AI-VISION1 Memory update

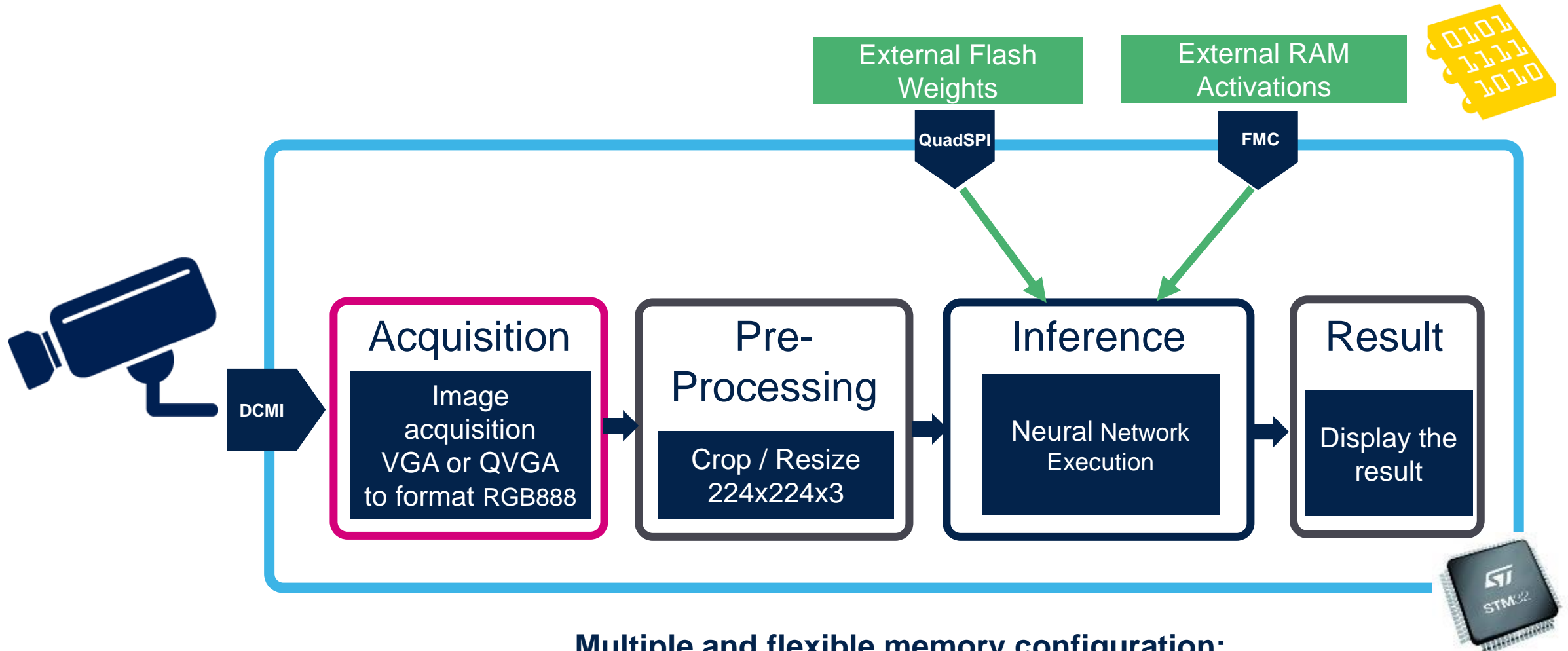


Optimized quantized model is fitting in internal memory: 153 ms

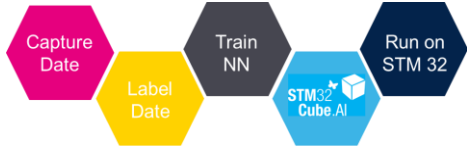


# FP-AI-VISION1

## Memory flexibility



**Multiple and flexible memory configuration:  
full internal, mixed, copy-before-run...**



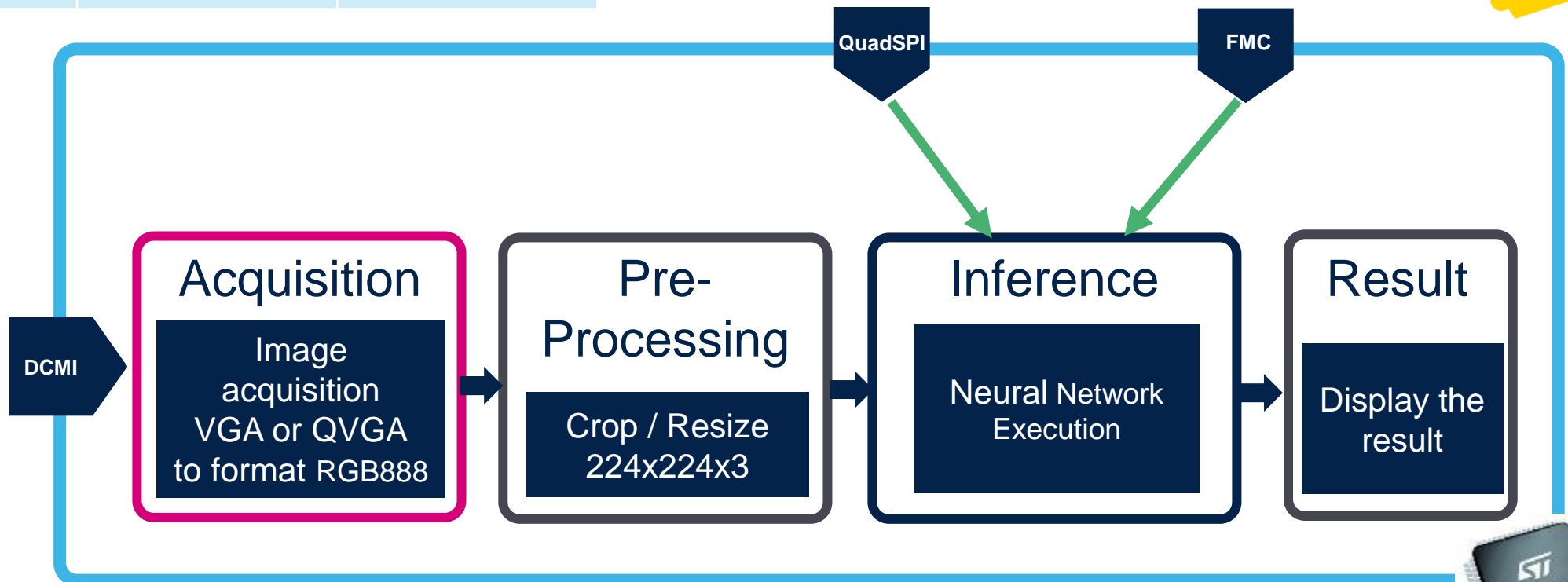
# FP-AI-VISION1

## Memory flexibility

| Memory    | Full internal | Full external |
|-----------|---------------|---------------|
| Inference | 155 ms        | 178 ms        |

External Flash Weights

External RAM Activations



Values measured on FP-AI-VISION1 1.0.0 with the optimized MobileNet Derivative model  
 Inference times on the STM32H7 @ 400MHz, Dual-QSPI DDR @ 50MHz, SDRAM @ 100MHz

# Demo

STM32    
Cube.AI





## Handwriting character recognition

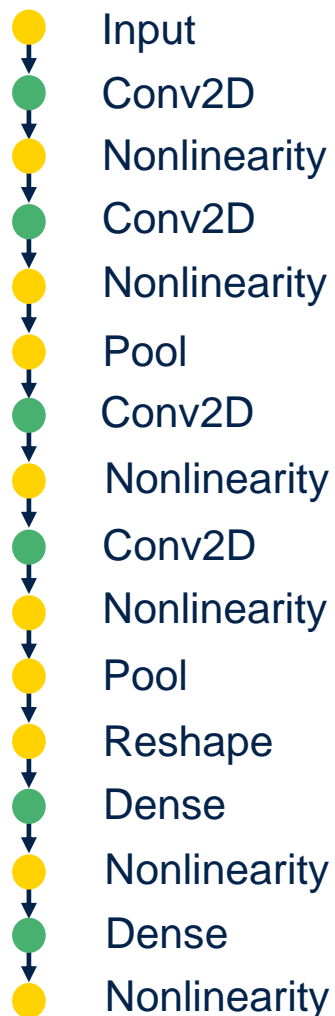


### Neural Network

- ST CNN
- EMNIST dataset (36 classes)

### Implementation

- Exploits touch screen captured as image of size 32x32
- 36 classes: numbers and capital letters



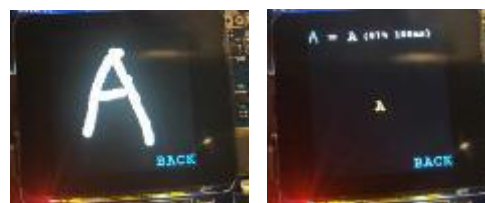
### STM32 Cube.AI NN

- Computational complexity 73k MACC
- Memory footprint: 26 KB RAM, 291 KB Flash



### Performance on STM32L562

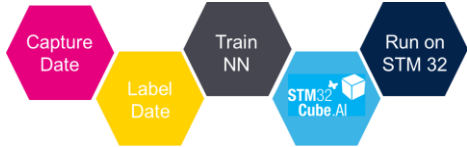
- 1 inference per image
- STM32L562 110MHz Cortex-M33



# Conclusion

STM32    
Cube.AI





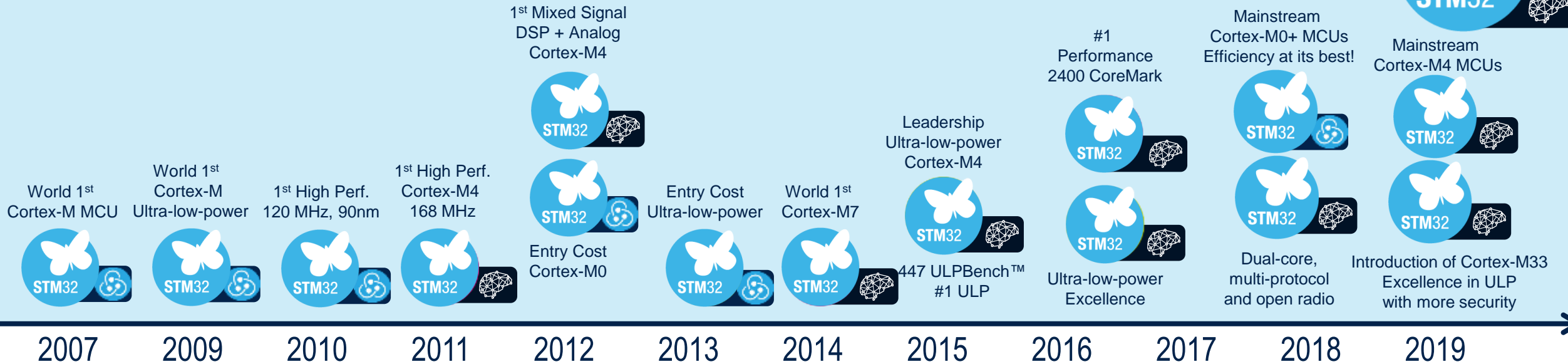
# Making AI accessible now

## Leader in Arm® Cortex®-M 32-bit General Purpose MCU

Compatible with **Deep Learning**  
STM32Cube.AI ecosystem



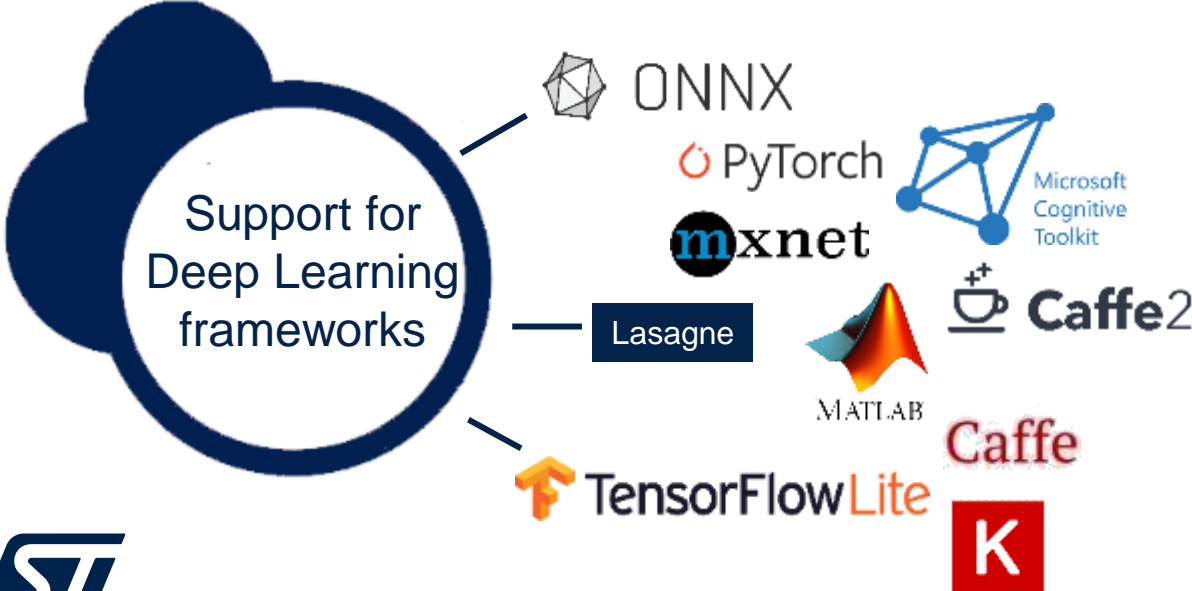
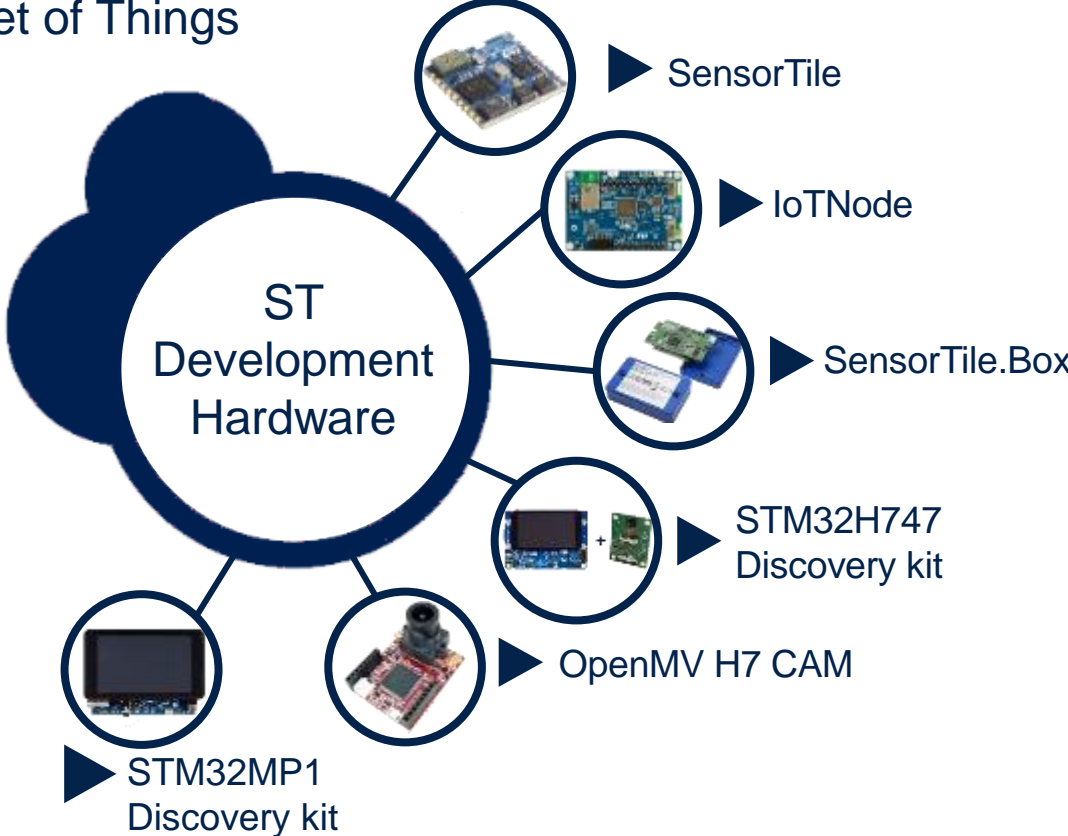
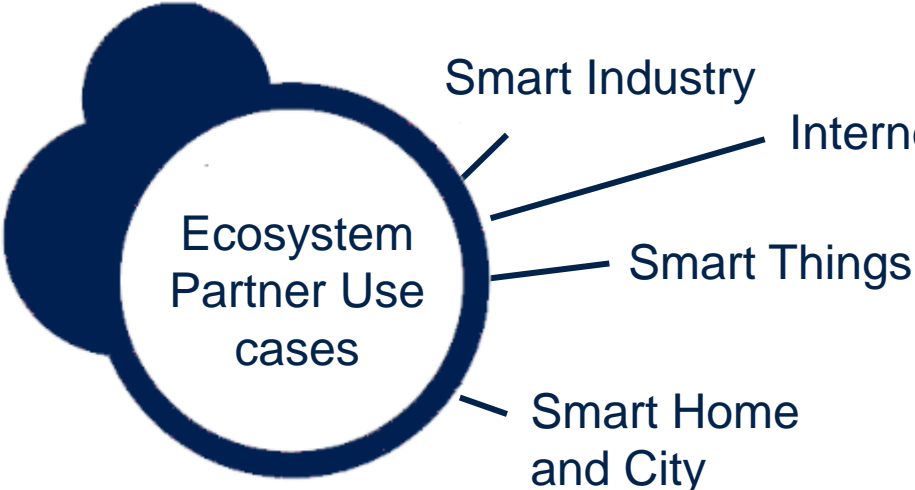
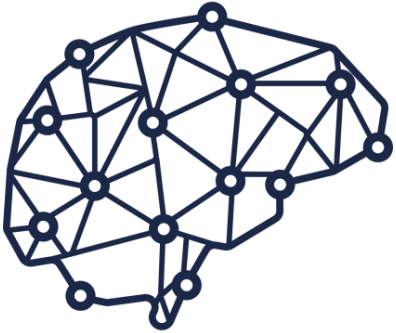
Compatible with **Machine Learning**  
Partner ecosystems

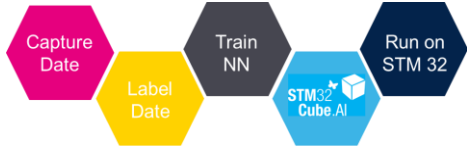


More than 40,000 customers

Over 4 Billion STM32 shipped since 2007

# Access the STM32Cube.AI ecosystem





# ST toolbox for Neural Networks

## More than just a conversion tool



- STM32Cube.AI to convert an NN in **optimized code**
- **Interoperability** with state-of-the-art Deep Learning design frameworks
- Support of **quantization**
  - Decrease memory requirement up to ↘÷4
  - Decrease latency and power consumption up to ↘÷3



- Software examples for Quick prototyping
- **Audio, Motion and Vision** Function packs on **ST development hardware**
- **STM32 Community with dedicated Neural Networks topic**

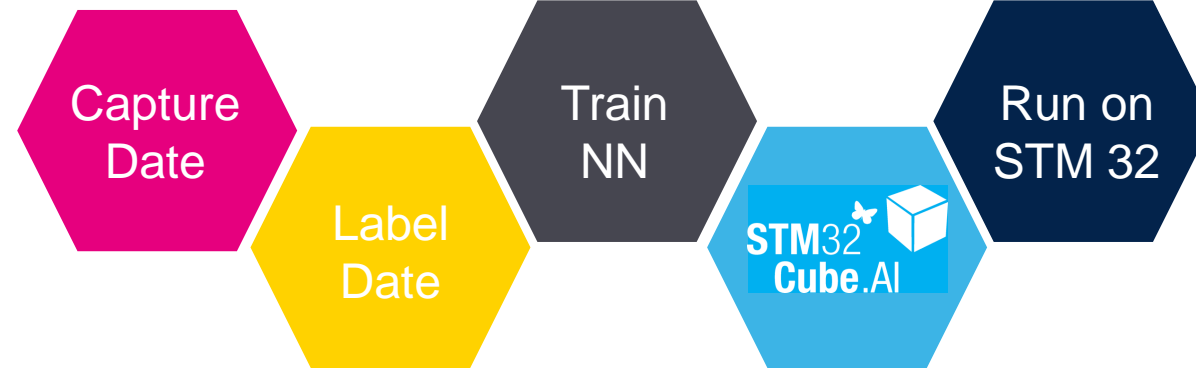


- **AI Partner Program**
- Expertise in **Machine Learning and STM32 solutions**

# For more information



[www.st.com/STM32CubeAI](http://www.st.com/STM32CubeAI)





# Thank you

© STMicroelectronics - All rights reserved.

ST logo is a trademark or a registered trademark of STMicroelectronics International NV or its affiliates in the EU and/or other countries.

For additional information about ST trademarks, please refer to [www.st.com/trademarks](http://www.st.com/trademarks).

All other product or service names are the property of their respective owners.



life.augmented