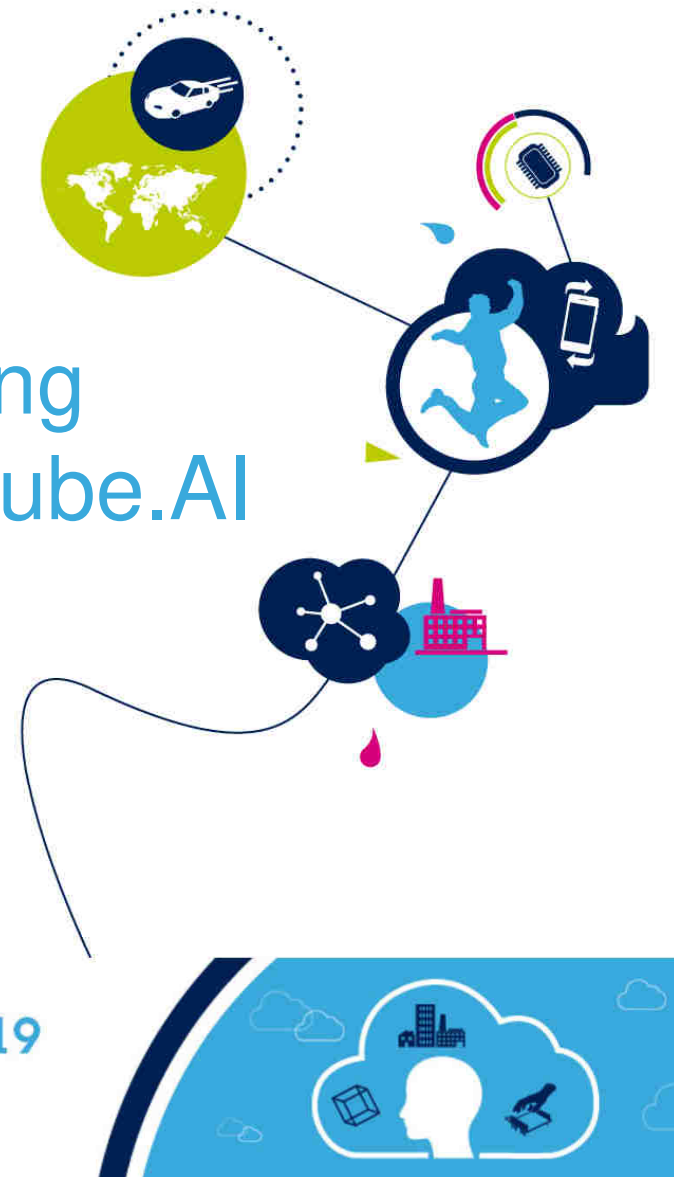


Artificial Neural Network Mapping Made Simple with the STM32Cube.AI

Markus Mayr
Product Marketing Manager, MCU



Technology Tour 2019
Minneapolis, MN | October 24



Artificial Intelligence (AI)

2

- AI is a superset of all the studies where machines mimic cognitive “human” capabilities. For example:
 - Interaction with the environment
 - Knowledge representation
 - Perception
 - Learning
 - Computer vision
 - Speech recognition
 - Problem solving and many more.
- Main ingredients
 - Computer science
 - Statistics
 - Mathematics



Artificial Intelligence (AI)

3

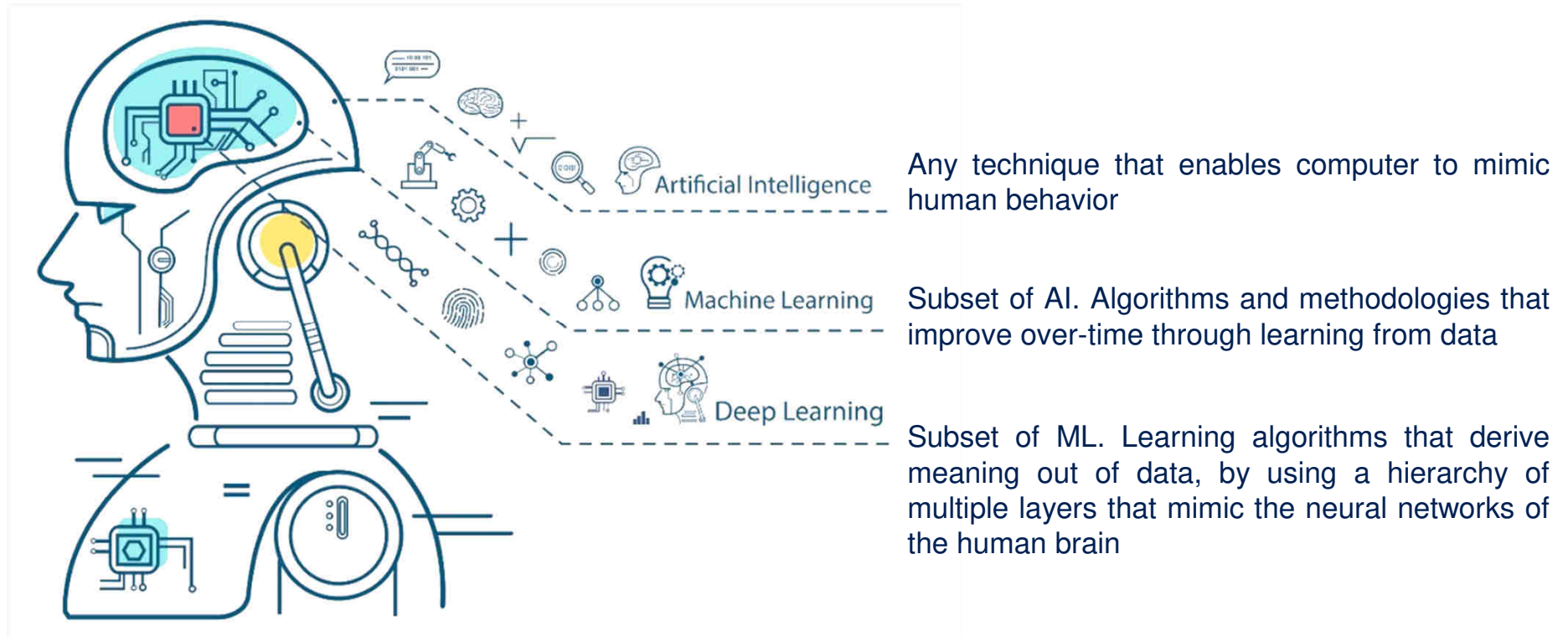
- **Main use cases in our everyday life:**

- Face/voice recognition
- Autonomous driving
- Stock market trading strategy
- Disease symptom detection
- Predictive maintenance
- Handwriting recognition
- Content distribution on social media
- Fraudulent credit card transaction
- Translation engines
- Shopping suggestions



Some definitions

4



Why Deep Learning is so Important

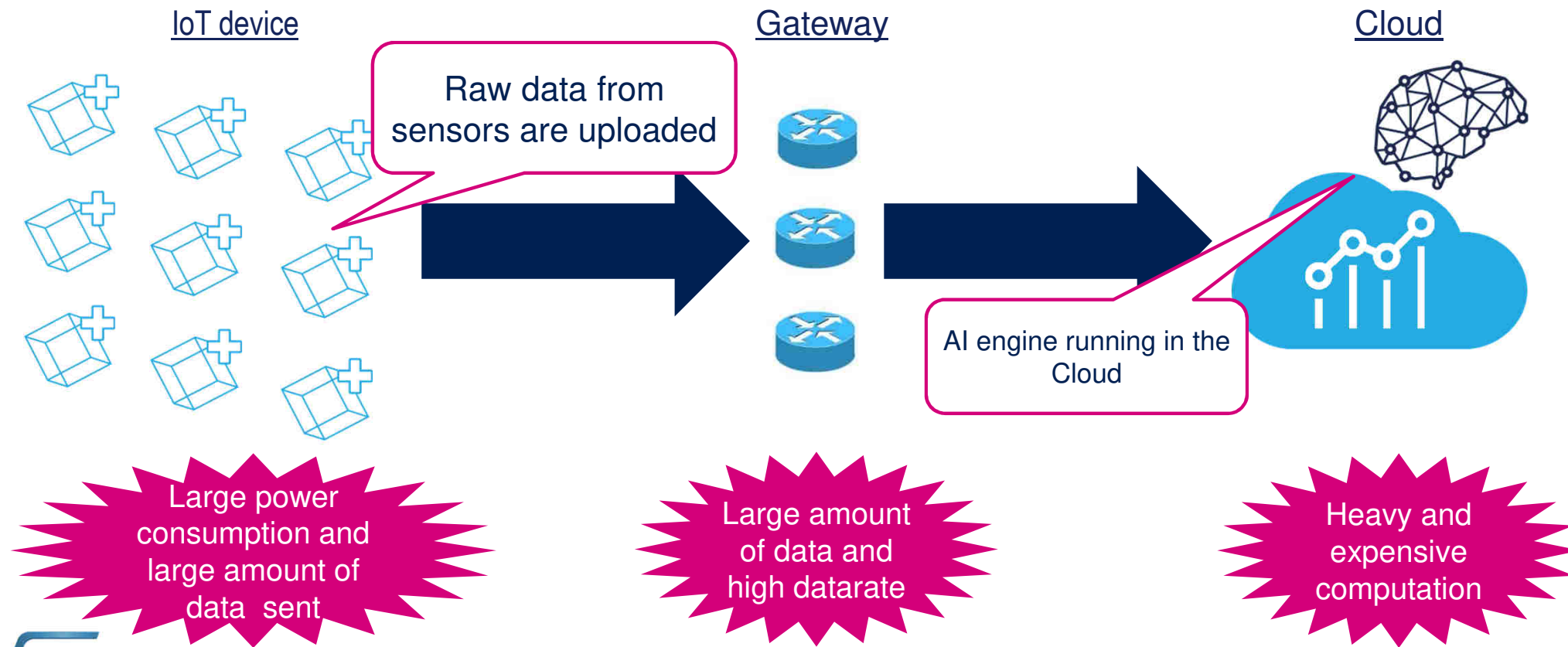
5

- Convolutional Deep Neural Networks outperform alternative methods on a number of tasks:

Problem	Dataset	Best Accuracy w/o CNN	Best Accuracy with CNN	Diff
Object classification	ILSVRC	73.8%	95.1%	+21.3%
Scene classification	SUN	37.5%	56%	+18.5%
Object detection	VOC 2007	34.3%	60.9%	+26.6%
Fine-grained class	200Birds	61.8%	75.7%	+13.9%
Attribute detection	H3D	69.1%	74.6%	+5.5%
Face recognition	LFW	96.3%	99.77%	+3.47%
Instance retrieval	UKB	89.3% (CDVS: 85.7%)	96.3%	+7.0%

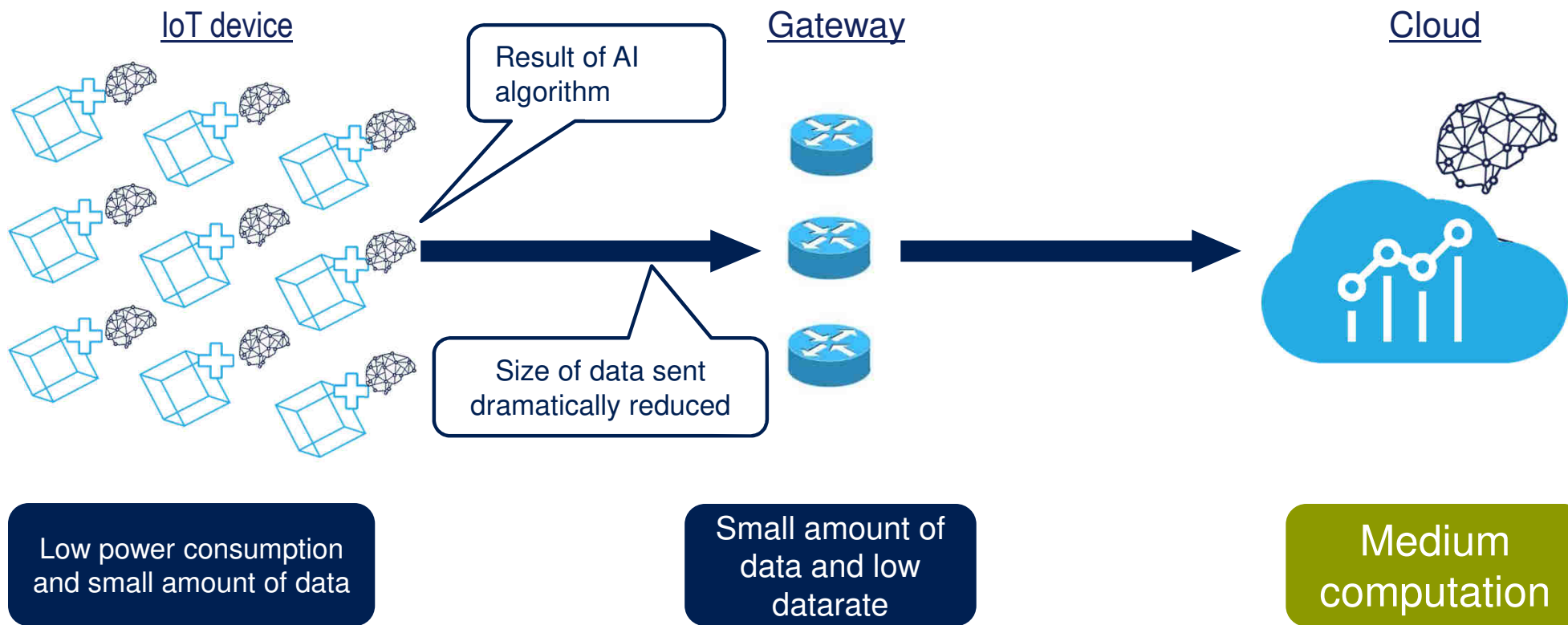
AI Cloud Computing

6



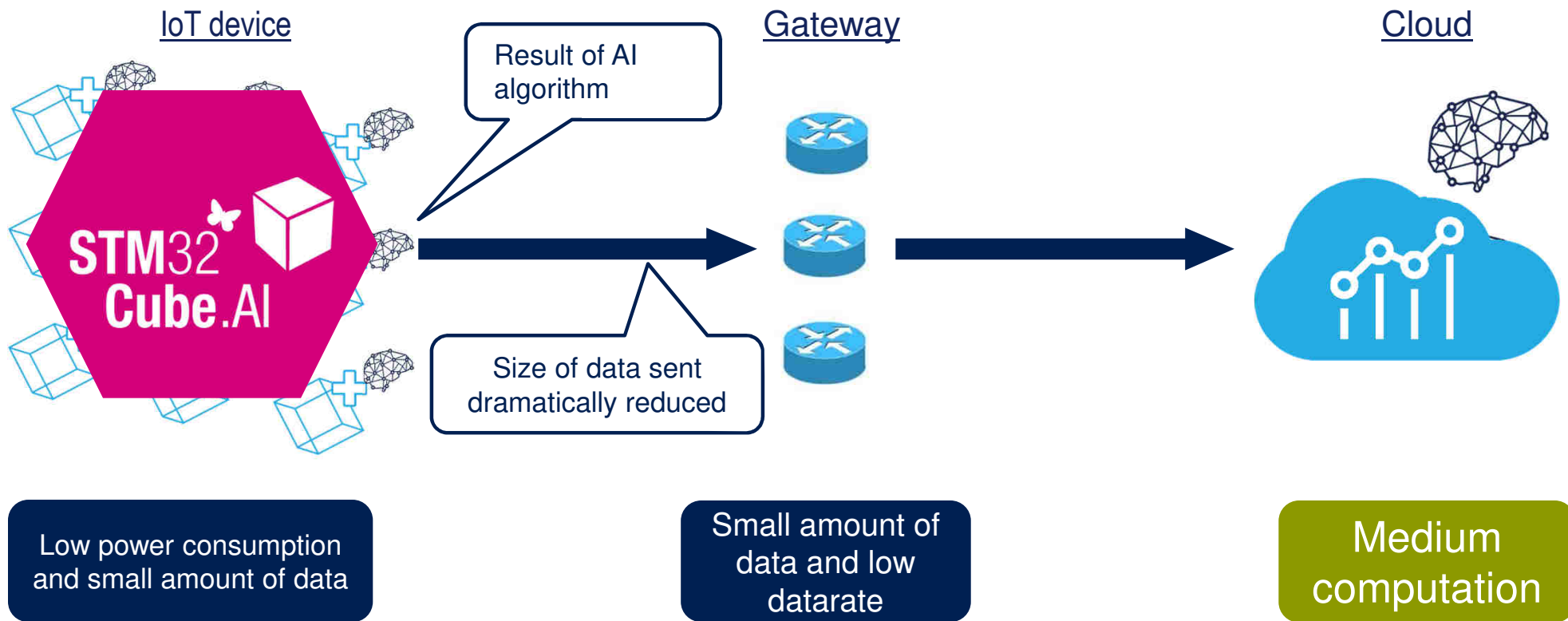
AI Edge Computing (Embedded)

7



AI Edge Computing (Embedded)

8



Distributed AI

9



High Bandwidth
High centralized computing power
Potentially high latency



Reduced bandwidth
Lower centralized computing power
Real-time response
Preserved Privacy

Neural Networks on STM32

Simple, fast, optimized



 **STM32**
Cube.AI 



The Key Steps Behind Neural Networks

11



Neural Network (NN) Model Creation



Operating Mode

Capture data



1

2



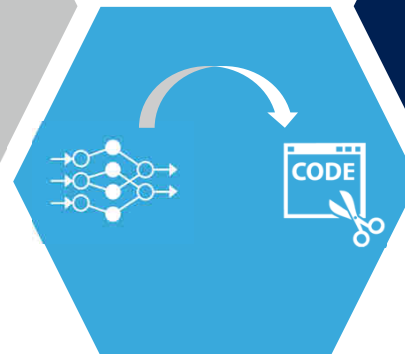
Clean, label Data
Build NN topology

Train NN Model



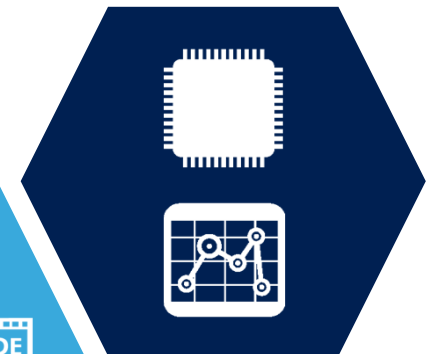
3

4



Convert NN into
optimized code for MCU

Process & analyze
new data using trained NN

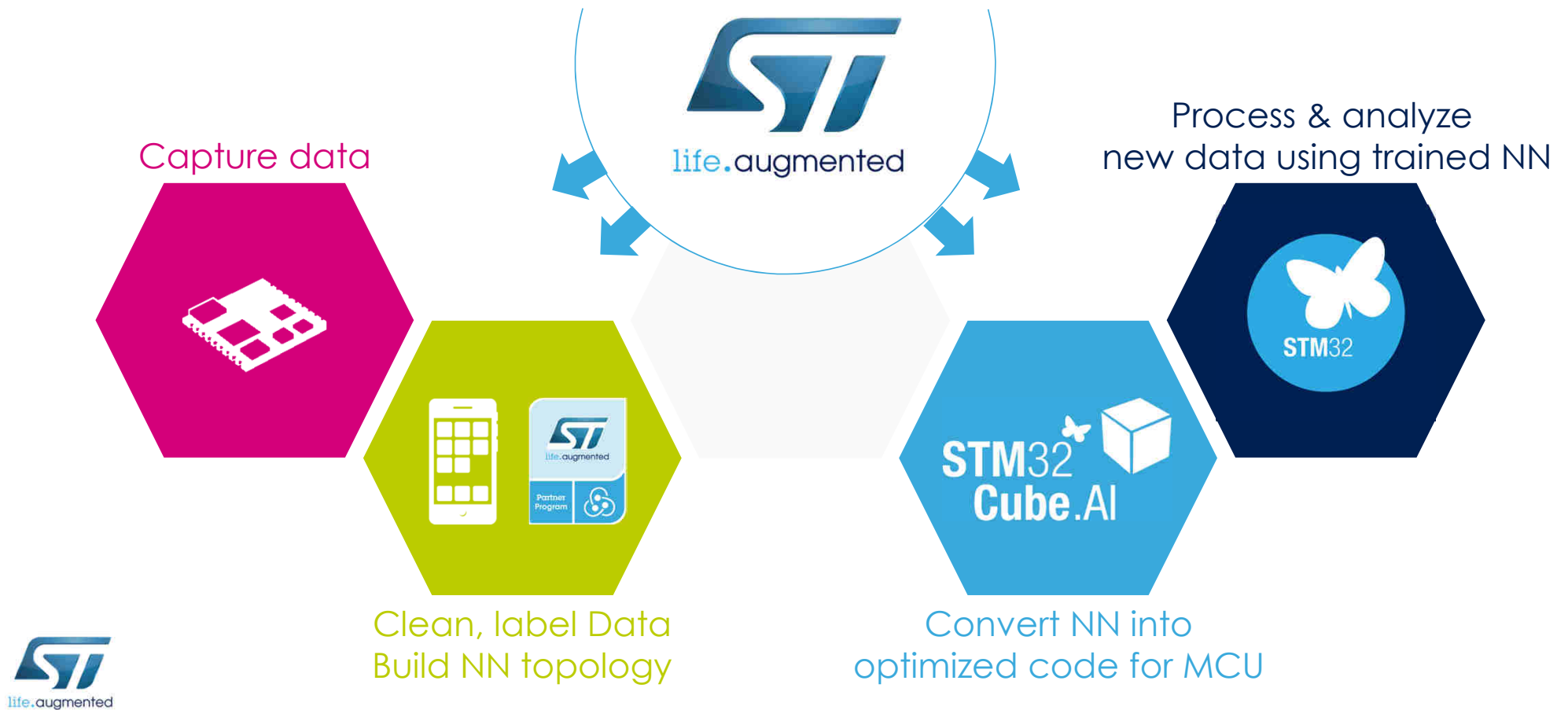


5



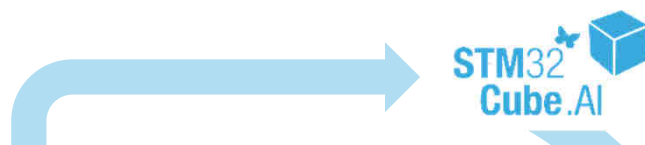
ST Toolbox for Neural Networks

12



STM32 Solutions for Embedded AI

Extensive toolbox to easily create your AI application



STM32
Cube.AI

AI extension for STM32CubeMX
to **map pre-trained Neural Networks**

Neural Networks on STM32
Simple, fast, optimized



STM32
Cube.AI



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












STM32 **Community** with dedicated
Neural Network topic



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

STM32 AI Typical Applications

14

Low	Medium	High
   <ul style="list-style-type: none"> • Sensor analysis • Activity recognition (motion sensors) • Stress analysis or attention analysis 	   <ul style="list-style-type: none"> • Audio & sound • Speech Recognition • Object detection 	   <ul style="list-style-type: none"> • Object detection / classification / tracking • Natural Language Understanding / Speech Synthesis

10s MOPs

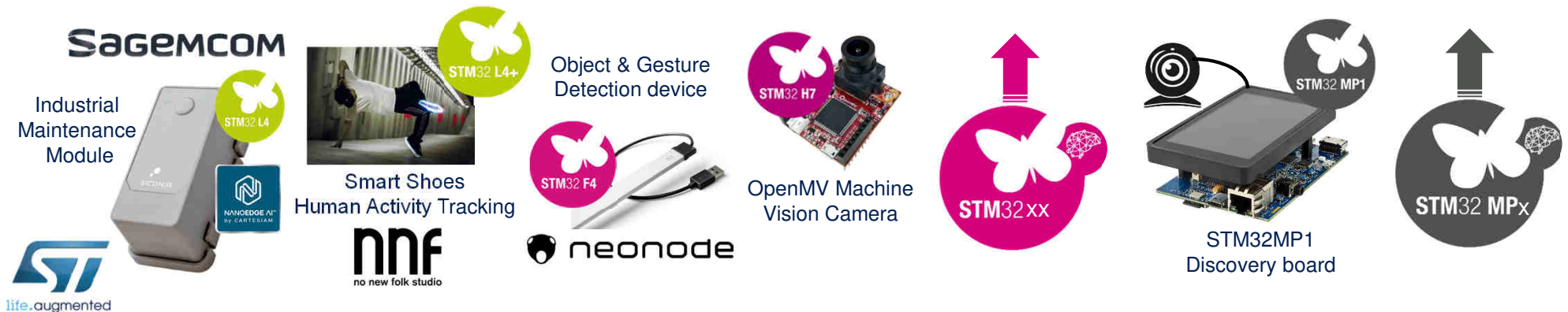
GOPs

0.5-1 TOPs

1-2 TOPs

MCU

From IP embedded in MCU/MPU to dedicated SOC





STM32CubeMX Extension

AI Conversion Tool

15

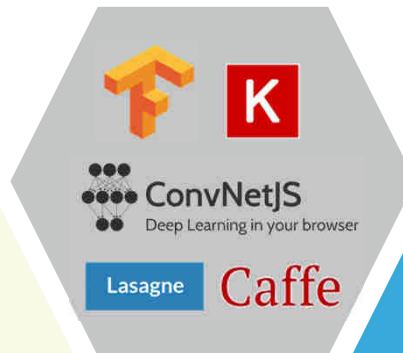
Input your framework-dependent, pre-trained Neural Network into the **STM32Cube.AI** conversion tool

Automatic, fast generation of an STM32-optimized library

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



Train NN Model



* TensorFlow used as a Keras backend.
Not all operators accessible to MCUs

Process & analyze new data using trained NN



STM32
Cube.AI

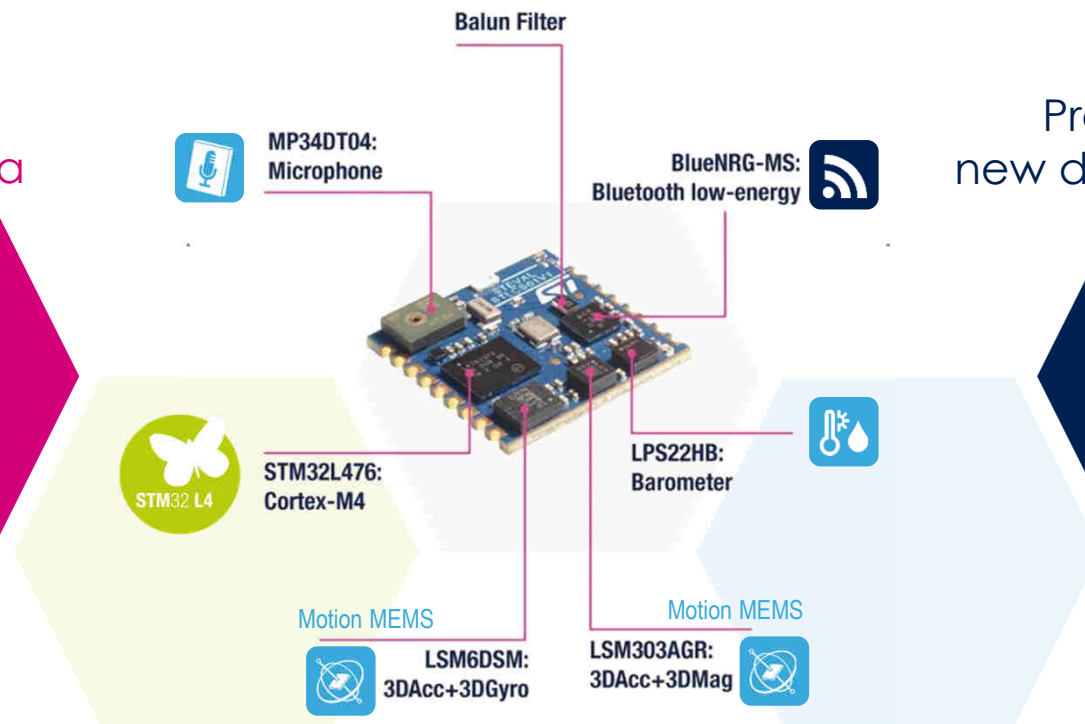
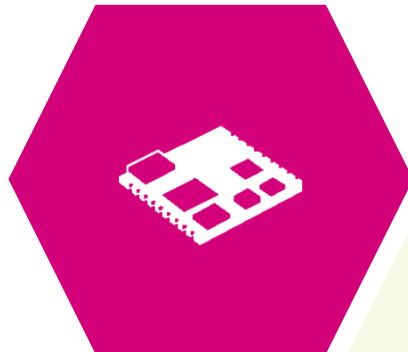
Convert NN into optimized code for MCU



Form Factor Hardware to Capture and Process Data

16

Capture data



Process & analyze
new data using trained NN



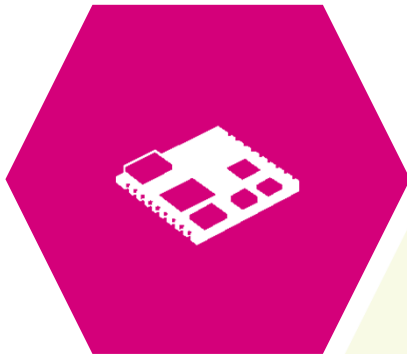


Form Factor Hardware

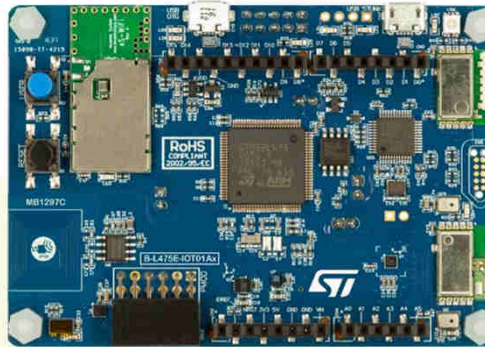
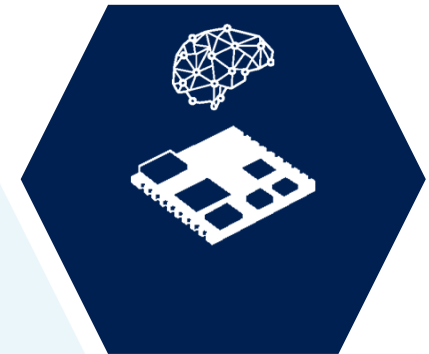
AI IoT Node for More Connectivity

17

Capture data



Process & analyze
new data using trained NN



More debug capabilities

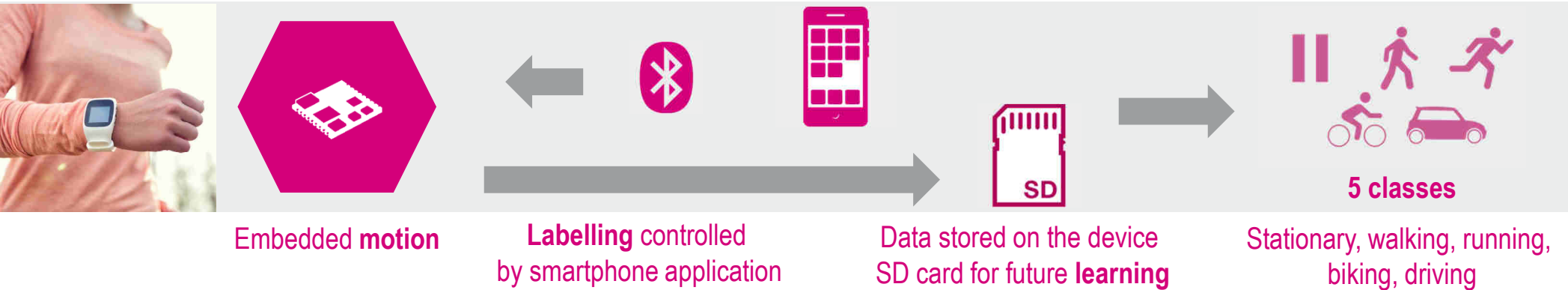
- Integrated ST-Link/V2.1
- PMOD extension connector
- Arduino Uno extension connectors



Human Activity Recognition (HAR)

Motion Example in FP-AI-SENSING1 Package

18





Human Activity Recognition

IGN (5 classes)

19

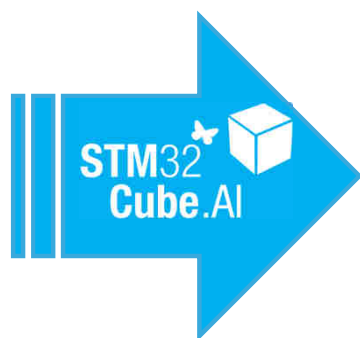
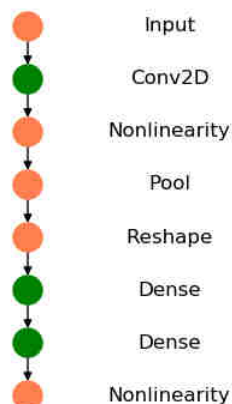
Use Case #2: HAR Human Activity Recognition Ignatov on SensorTile

Neural Network

- Derived from a published paper Keras model
- ST proprietary dataset of 2.4M samples

Implementation

- Exploits 3-axis accelerometer data
- 5 classes: stationary, walking, running, biking, driving
- Pre/Post-processing: filtering gravity, reference rotation, temporal filter



STM32 Cube.AI NN

- Computational complexity 14k MACC
- Memory footprint: 1.8 KB RAM, 12 KB Flash



Performance on Sensor Tile

- STM32L476 80MHz Cortex-M4F
- Use case: 1 classification/sec
- Pre/Post-processing: 0.02 MHz NN processing: 0.35 MHz
- Power consumption (1.8 V)
 - System: 580 uA (with optim BLE)
 - STM32: 510 uA



Audio Scene Classification (ASC)

Audio Example in FP-AI-SENSING1 Package

20

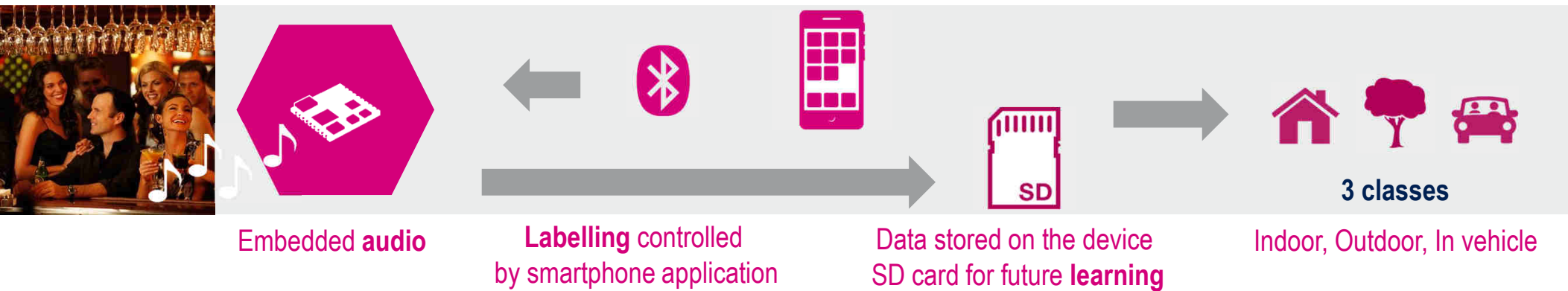




Image Classification

21

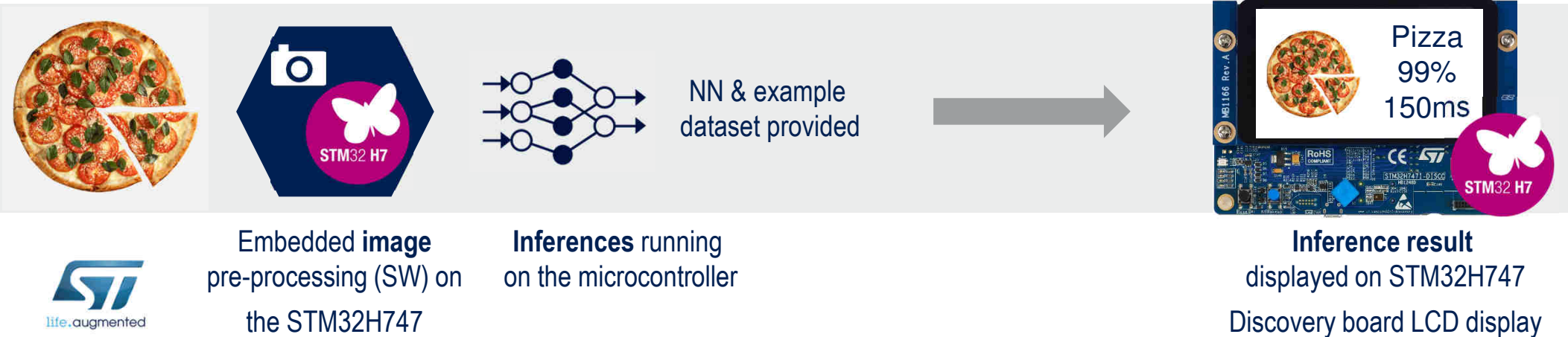
Vision Example in FP-AI-VISION1 Package

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





Food Recognition

22

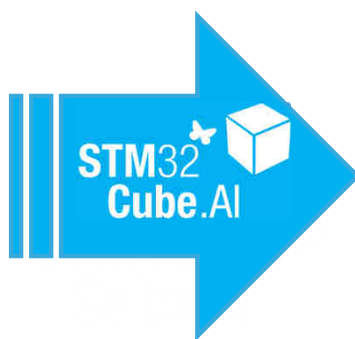
Fast Downsampling MobileNet Food Recognition on STM32H747 Dual-Core Discovery board

Neural Network

- FD-MobileNet topology from public paper applied to food
- Dataset

Implementation Details

- Uses Camera either in continuous or one shot mode
- Floating Point or mixed model Floating/Fix Point
- 18 food classes



STM32 Cube.AI NN

- Memory footprint: 205 KB RAM, 191 KB Flash

Performance on STM32H747

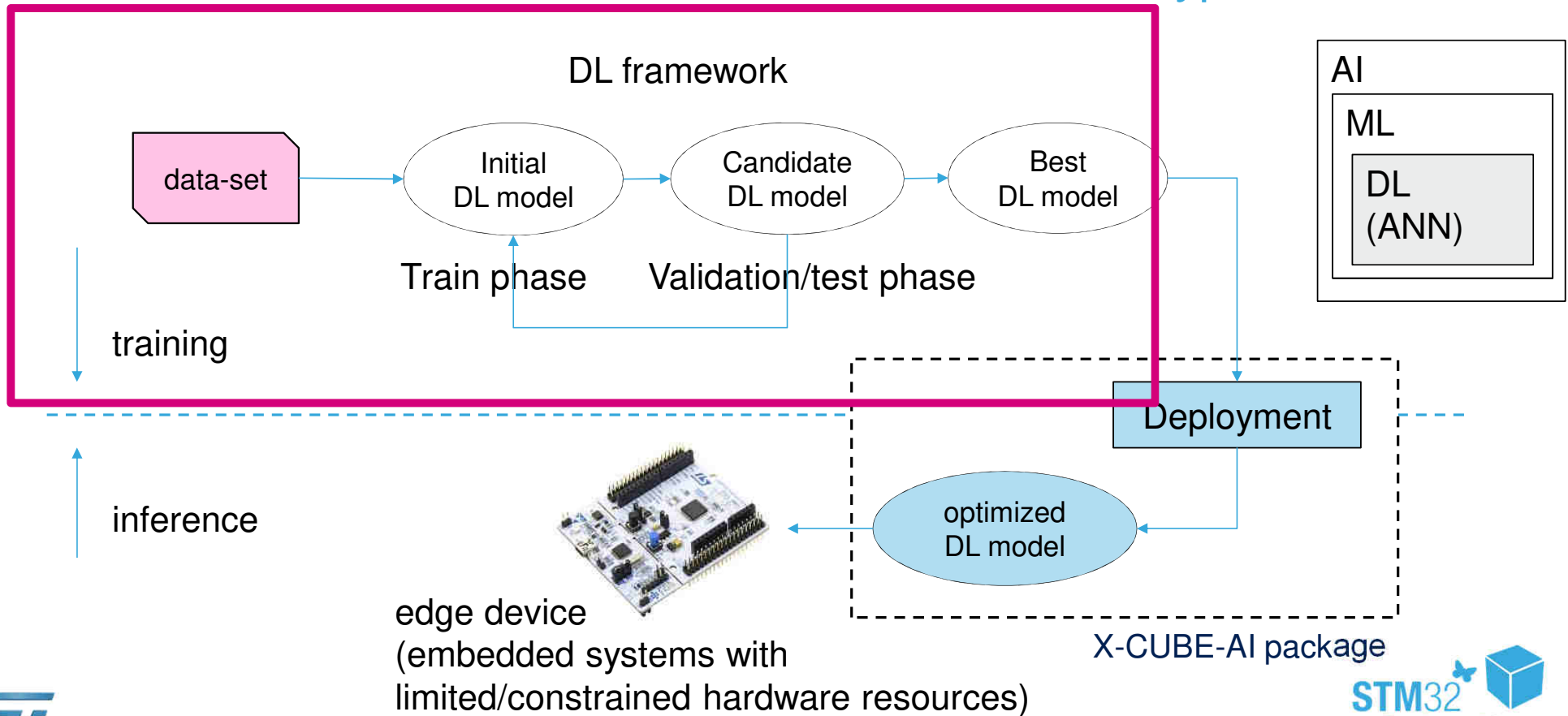
- 1 inference per image
- STM32H747 400 MHz Cortex-M7F
- Mix model Fix/Floating Point
 - 6.2 MHz / 150 ms per inference
 - Accuracy: 78.8%



X-CUBE-AI Positioning

in a typical DL flow

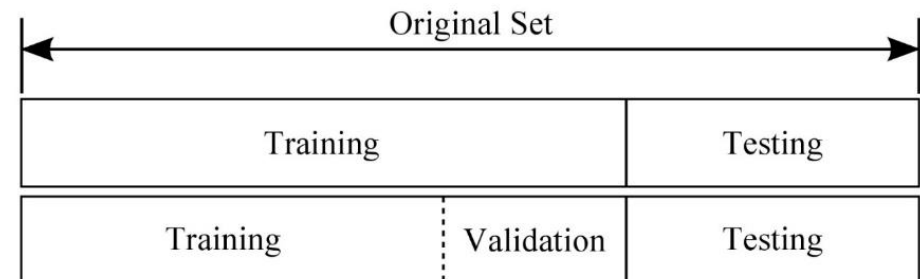
23



Learning datasets

24

- One of the difficulties of working with NN is to have a big dataset. In general 90% of the time spent for a Deep Learning project is related to the creation of the Dataset.
- For Image Processing and few other applications there are already available big databases (with labels), both free or paid basis. On other hand for new and “niche” applications the designer needs to create the database, that becomes an high value asset.
- The dataset is then divided into
 - Training set
 - Training set
 - Validation set
 - Testing set



Data-driven approach

25

- We provide many examples of each class to the computer and then develop **machine learning algorithms** that look at these examples and **learn about the visual appearance** of each class.

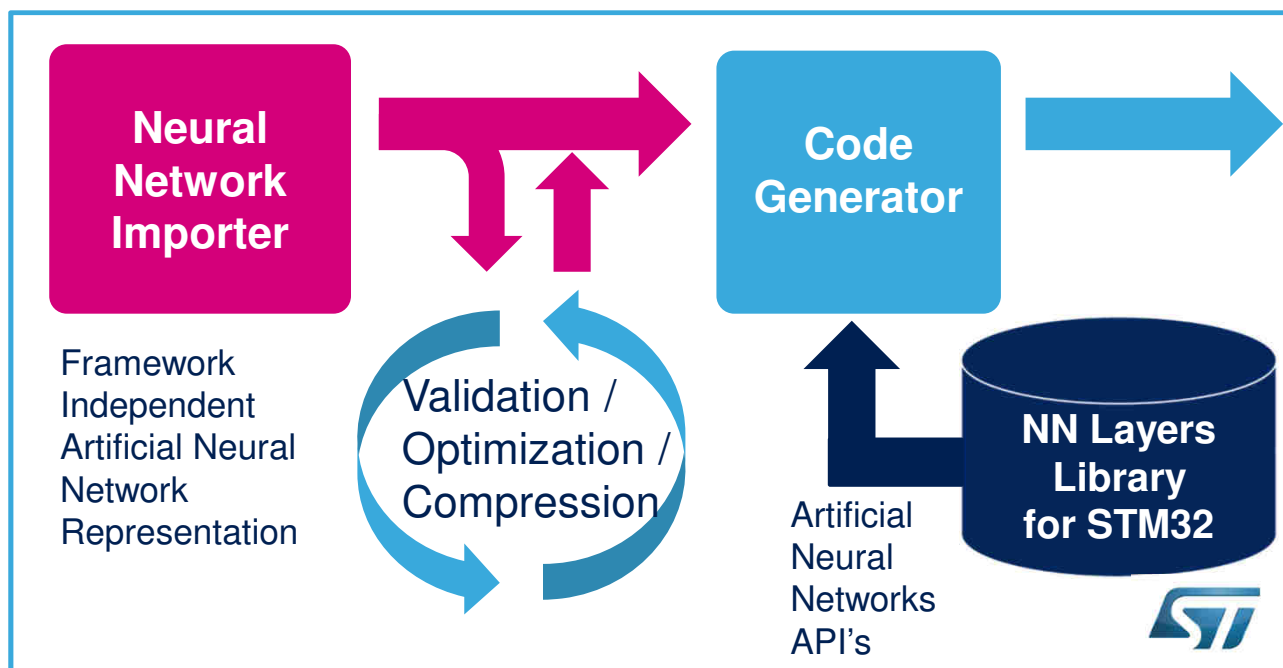


X-Cube-AI : Architecture

26

Off-the-shelf :
Pre-trained Artificial
Neural Network Model

Deep Learning
Framework dependent



Embedded Solution
Optimized Artificial
Neural Network Code
generated for STM32



This optimized STM32 Artificial neural network model can be included into the user project (using KEIL, IAR, OpenSTM32) and can be compiled and ported onto the final device for field trials

MHz and embedding a floating point unit (FPU). The family incorporates high-speed embedded memories (up to 64 Kbyte of Flash

Graphic Summary AI Summary

K Keras

Minimum Ram: 196 Bytes
Minimum Flash: 15.20 KBytes

C:\Users\ledonger\Documents\deepnet_relu.h5

MCUs List: 627 items

Display similar items

	Part No	Refere	Marketing	Unit Price for 1...	Board	Package	Flash	RAM	IO	Freq	GFX S	HMAC	MD5	SH
☆	STM32F301C6	STM3...	Active	1.596		LQFP48	32 kByt...	16 kBytes	37	72 MHz	0.0	0	0	0
☆	STM32F301C8	STM3...	Active	1.666		LQFP48	64 kByt...	16 kBytes	37	72 MHz	0.0	0	0	0
☆	STM32F301K6	STM3...	Active	1.272		WLCSP49	64 kByt...	16 kBytes	37	72 MHz	0.0	0	0	0
☆	STM32F301K6	STM3...	Active	1.272		LQFP32	32 kByt...	16 kBytes	25	72 MHz	0.0	0	0	0
☆	STM32F301K8	STM3...	Active	1.342		UFQFPN32	32 kByt...	16 kBytes	24	72 MHz	0.0	0	0	0
☆	STM32F301K8	STM3...	Active	1.342		LQFP32	64 kByt...	16 kBytes	25	72 MHz	0.0	0	0	0
☆	STM32F301R6	STM3...	Active	1.758		UFQFPN32	64 kByt...	16 kBytes	24	72 MHz	0.0	0	0	0
☆	STM32F301R6	STM3...	Active	1.758		LQFP64	32 kByt...	16 kBytes	51	72 MHz	0.0	0	0	0
☆	STM32F301R8	STM3...	Active	1.828		LQFP64	64 kByt...	16 kBytes	51	72 MHz	0.0	0	0	0
☆	STM32F302C6	STM3...	Active	1.712		LQFP48	32 kByt...	16 kBytes	37	72 MHz	0.0	0	0	0
☆	STM32F302C8	STM3...	Active	1.782		LQFP48	64 kByt...	16 kBytes	37	72 MHz	0.0	0	0	0
☆	STM32F302C8	STM3...	Active	1.782		WLCSP49	64 kByt...	16 kBytes	37	72 MHz	0.0	0	0	0
☆	STM32F302CB	STM3...	Active	1.99		LQFP48	128 kBy...	32 kBytes	37	72 MHz	0.0	0	0	0
☆	STM32F302CC	STM3...	Active	2.288		LQFP48	256 kBy...	40 kBytes	37	72 MHz	0.0	0	0	0

☐ Enable

Artificial Intelligence

☒ Enable

Model

Keras

Type

Saved model

Model

deepnet_relu.h5

Browse

Compression

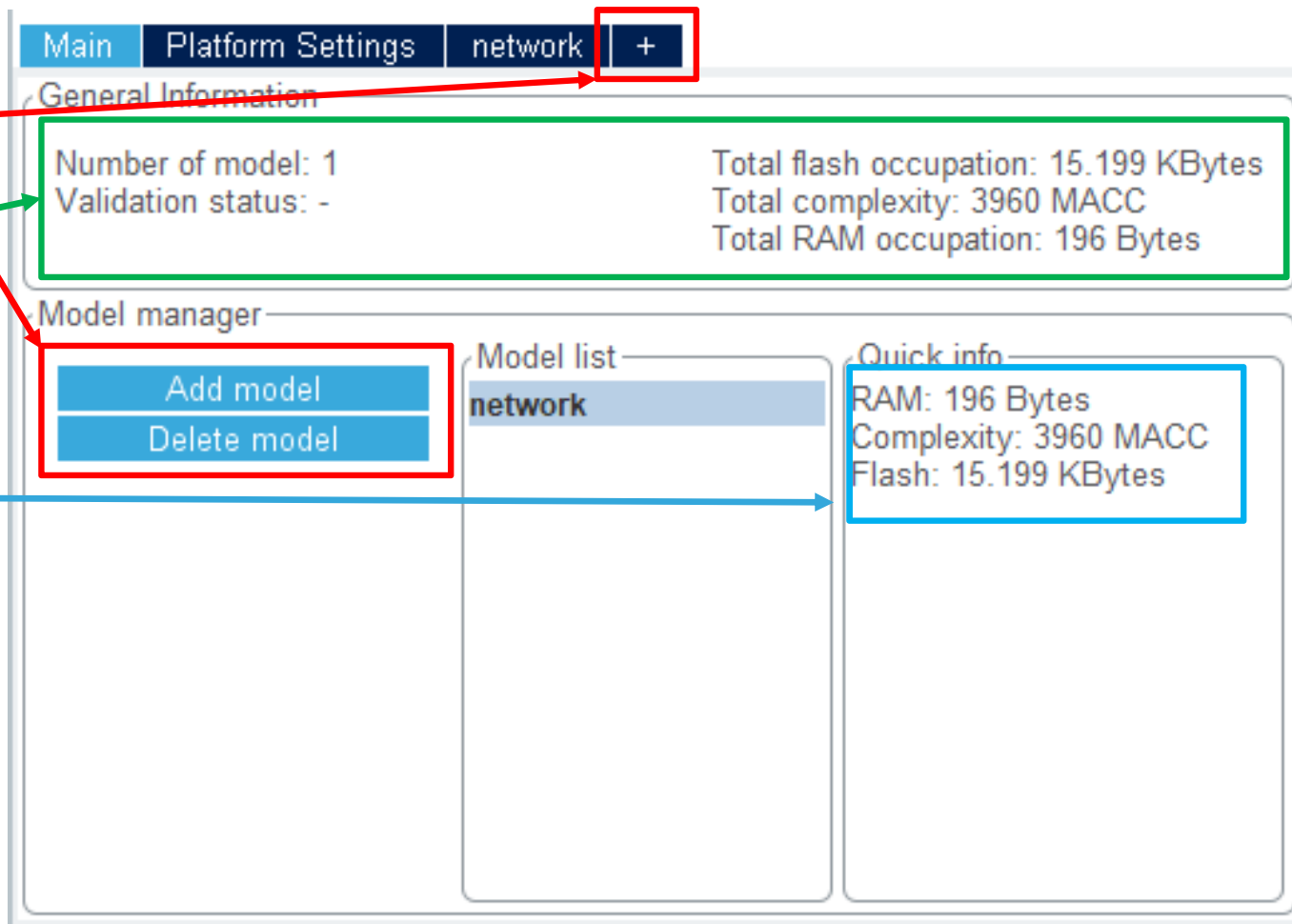
None

Analyze

Peripheral

ADC 12-bit	0	40
ADC 16-bit	0	36
AES	<input type="checkbox"/>	
CAN	0	3
COMP	0	7
CRYP	<input type="checkbox"/>	
DAC 12-bit	0	3

- Add/Delete models
- Get general information
- Have a quick look at different models



The screenshot shows the X-Cube-AI Main tab interface. The top navigation bar includes 'Main', 'Platform Settings', 'network', and a '+' button. The 'General Information' section displays model statistics. The 'Model manager' section contains 'Add model' and 'Delete model' buttons. The 'Model list' shows the 'network' model. The 'Quick info' section provides a summary of resources.

General Information	
Number of model: 1	Total flash occupation: 15.199 KBytes
Validation status: -	Total complexity: 3960 MACC
	Total RAM occupation: 196 Bytes

Model manager	
<div>Add model</div> <div>Delete model</div>	<div>Model list</div> <div>network</div>

Quick info
RAM: 196 Bytes Complexity: 3960 MACC Flash: 15.199 KBytes

X-Cube-AI Detailed View

29

- Perform **analysis** to compute the model size, get an image of the network and the complexity
- Perform **validation on desktop**
- Perform **validation on target**
- Set a **compression** to reduce the model size (By reducing the accuracy of the model)

Main
Platform Settings
network
+

Model inputs
network

Keras
Saved model

Model: C:\Users\vedonger\Documents\deepnet_relu.h5
Browse

Browse

Command:

Validation status: Unknown
Complexity: -
Flash occupation: -
RAM: -
Actual compression: -

Compression: None

Show graph

Analyze

Validation from: Random numbers

Validate on desktop

Validate on target

Making AI Accessible Now

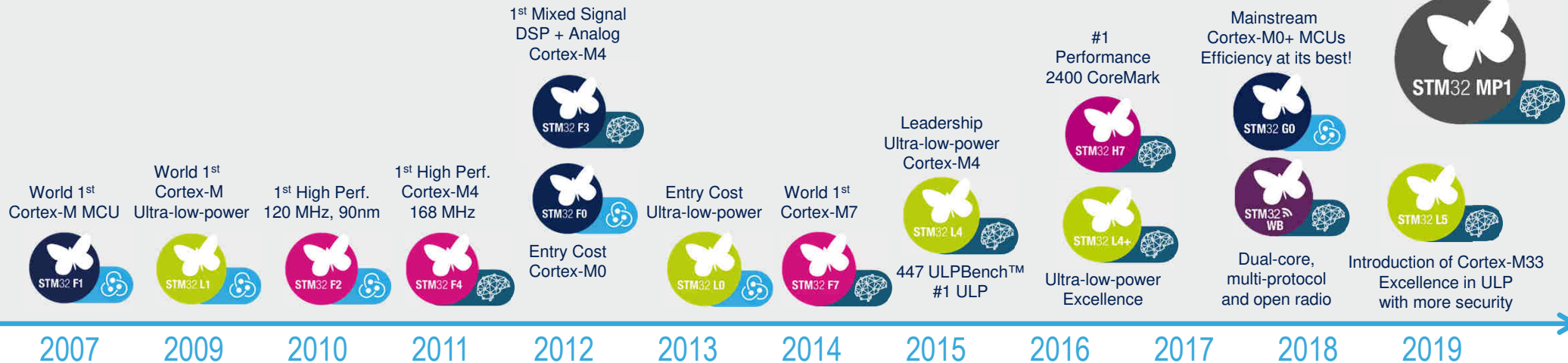
30

Leader in Arm® Cortex®-M 32-bit General Purpose MCUs

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



Compatible with **Machine Learning**
Partner ecosystems



More than 60,000 customers

Over 4 Billion STM32 shipped since 2007



STM32Cube.AI Roadmap

31

2019

Feb

Mar

Apr

May

Jun

Jul

Aug

Sept

Oct

Dec

2020

Jan

Feb

Mar

Apr

STM32
Cube.AI

- Market Introduction
- Floating Point Support

- Additional layers

- Fixed Point Quantization
- Command line interface
- UI Improvements
- Additional layers

- Integer Arithmetic Quantization
- Advanced External Flash Support

- ONNX introduction
- Additional layers
- Debug improvements

NN
Training
Tools
Supported



Lasagne

ConvNetJS
Deep Learning in your browser

Caffe

+ TensorFlow Lite

- Keras Fixed Point
- TensorFlow Lite Float

+ TensorFlow Lite

- TensorFlow Lite Integer arithmetic quantization



+

ONNX

moxnet

Microsoft
CNTK

Chainer

PyTorch



New layers & operators addition



ST
life.augmented



STM32 Solutions for AI

More Than Just the STM32Cube.AI

32

An extensive toolbox to support easy creation of your AI application

AI extension for STM32CubeMX

To map pre-trained Neural Networks onto the STM32

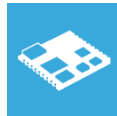


Function packs for Quick prototyping

Audio and motion examples

SensorTile reference hardware

To run inferences or data collection



... And more coming!



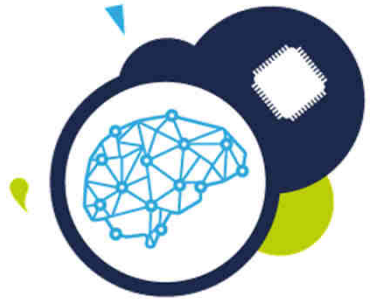
STM32 Community with dedicated Neural Networks topic

Mobile phone application

To collect and label data
To display the result of inference processing on the STM32



ST Partner Program with a dedicated group of Partners providing Neural Networks engineering services
Data scientists and Neural network architects



For More Information

33



www.st.com/STM32CubeAI

