

Ubiquitous AI

Smart and secure IoT as enabler contributors boosting servitization approaches

Claudio Marchisio

STMicroelectronics

*System Research & Applications /
AI SW & Tools*

Organizzato da



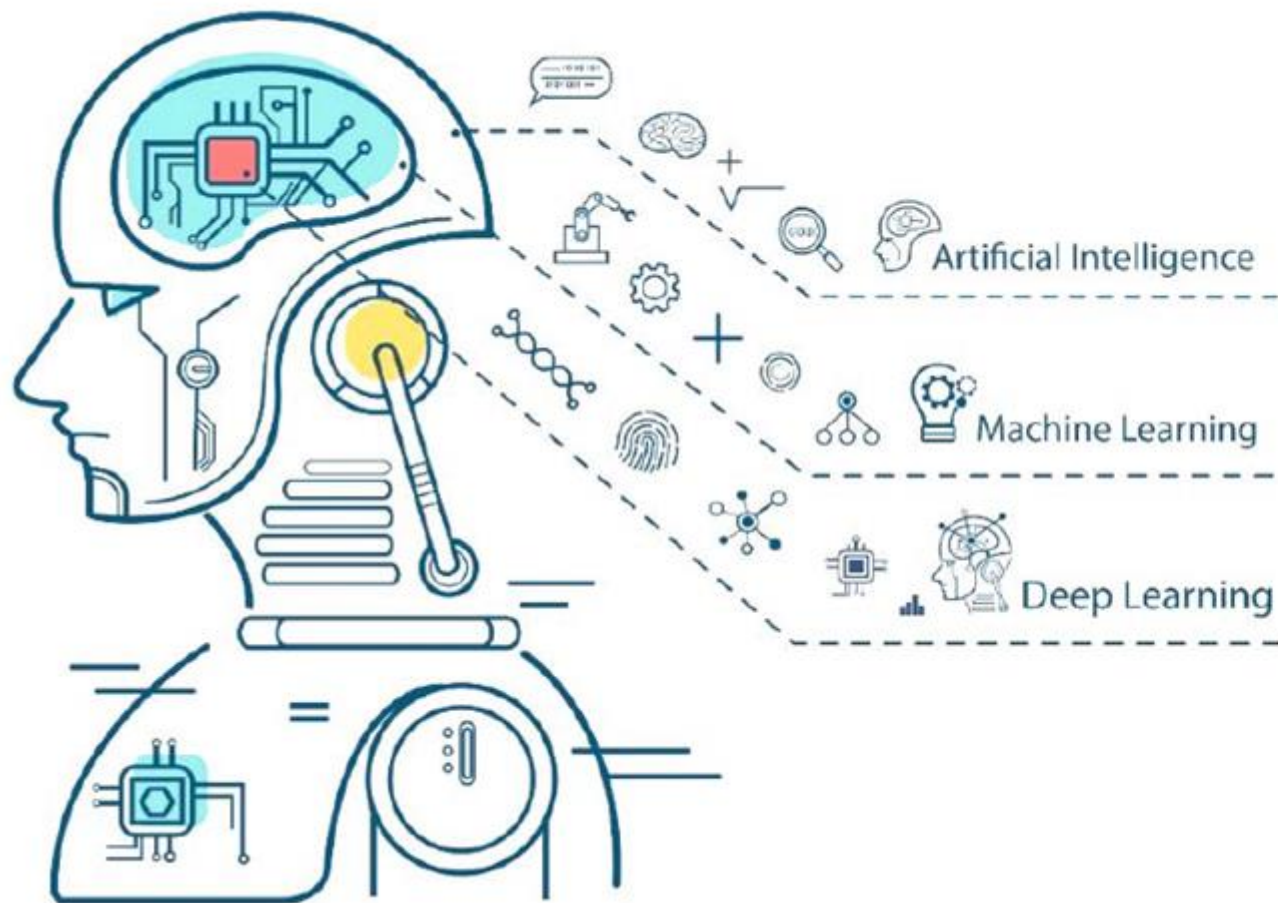
Edge AI



STM32 
Cube.AI















Artificial Intelligence

Some definitions



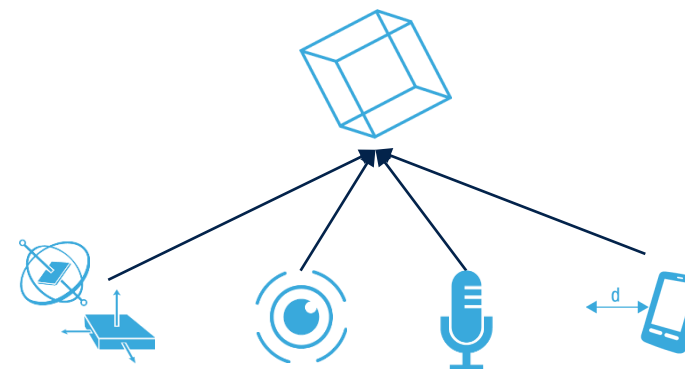
- AI: any technique which enables a computer to mimic human behavior
- ML is a subset of AI that provides systems the ability to automatically learn and improve from data without being explicitly programmed
- DL is a subset of ML, utilizes a hierarchical level of artificial neural networks to carry out the process of machine learning

The building blocks of the IoT

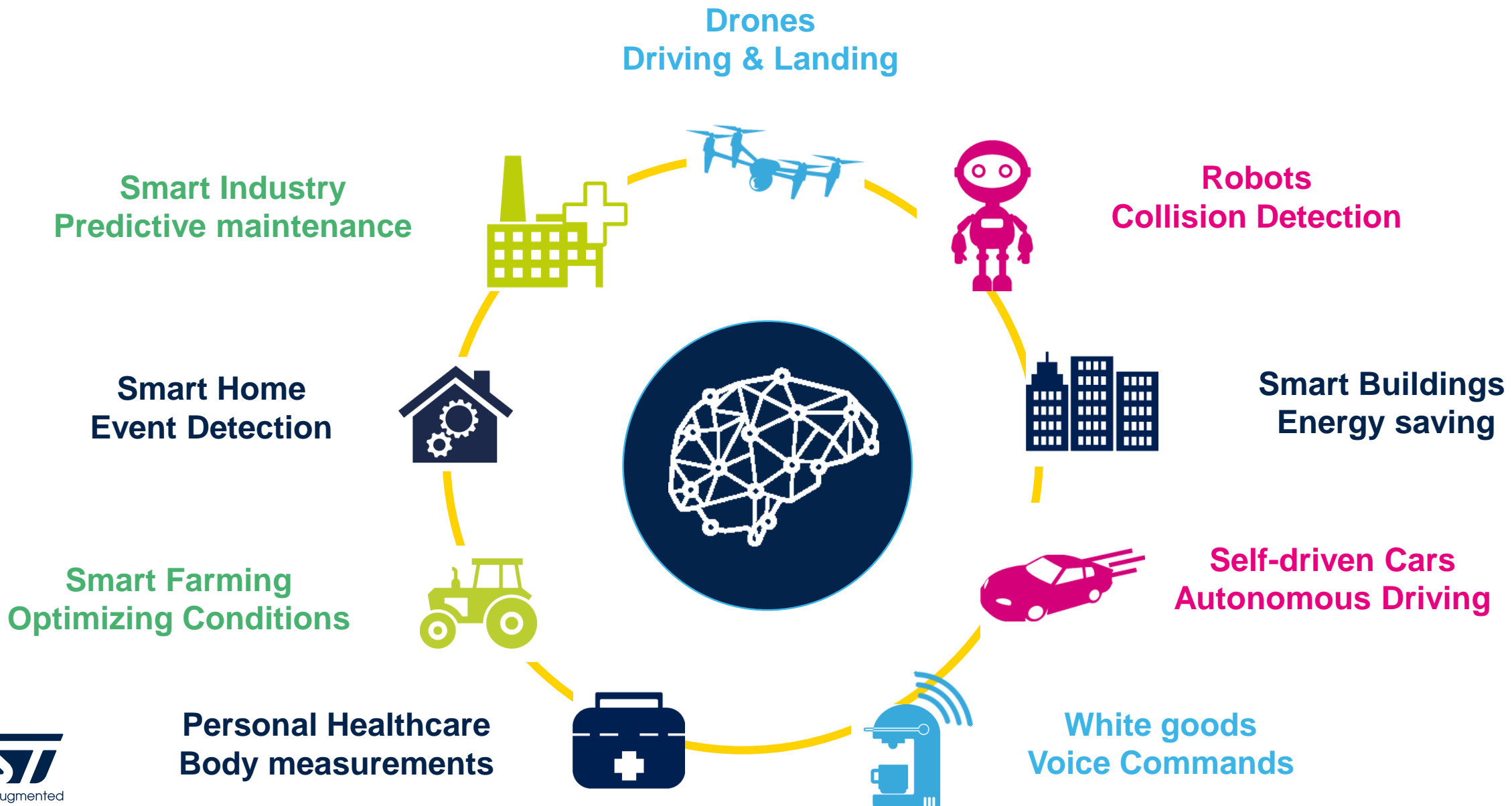
	Processing	Security	Sensing & Actuating	Connectivity	Conditioning & Protection	Motor Control	Power & Energy Management
Smart Things							
Smart Home & City	Ultra-Low Power to High Performance	Scalable security solutions	Full range of sensors and actuators	10 cm to 10 km	Nano Amps to Kilo Amps	Power conversion Monitoring Drivers	Nano Watt to Mega Watt
Smart Industry							

What is Edge AI?

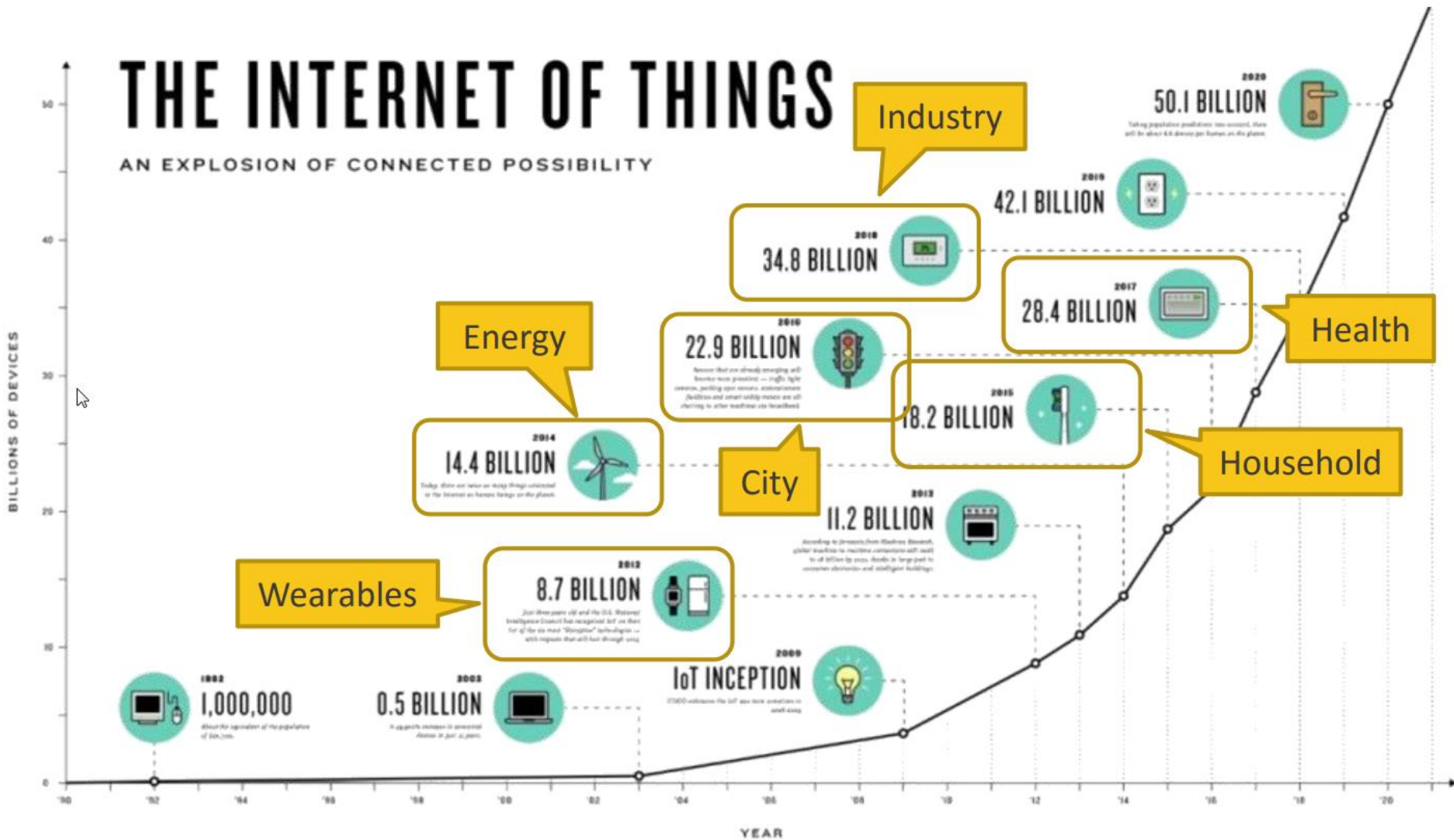
- **Edge AI** means that AI algorithms are executed locally on a hardware device to process the data generated by the attached sensors
- An **Edge AI device** processes data, extracts information and takes decisions without a connection
- Anyhow a connection (IoT nodes) is useful to:
 - receive FW updates, update Deep Learning models, etc.
 - send (reduced) data and results of local processing
 - receive commands



Edge AI means Ubiquitous AI

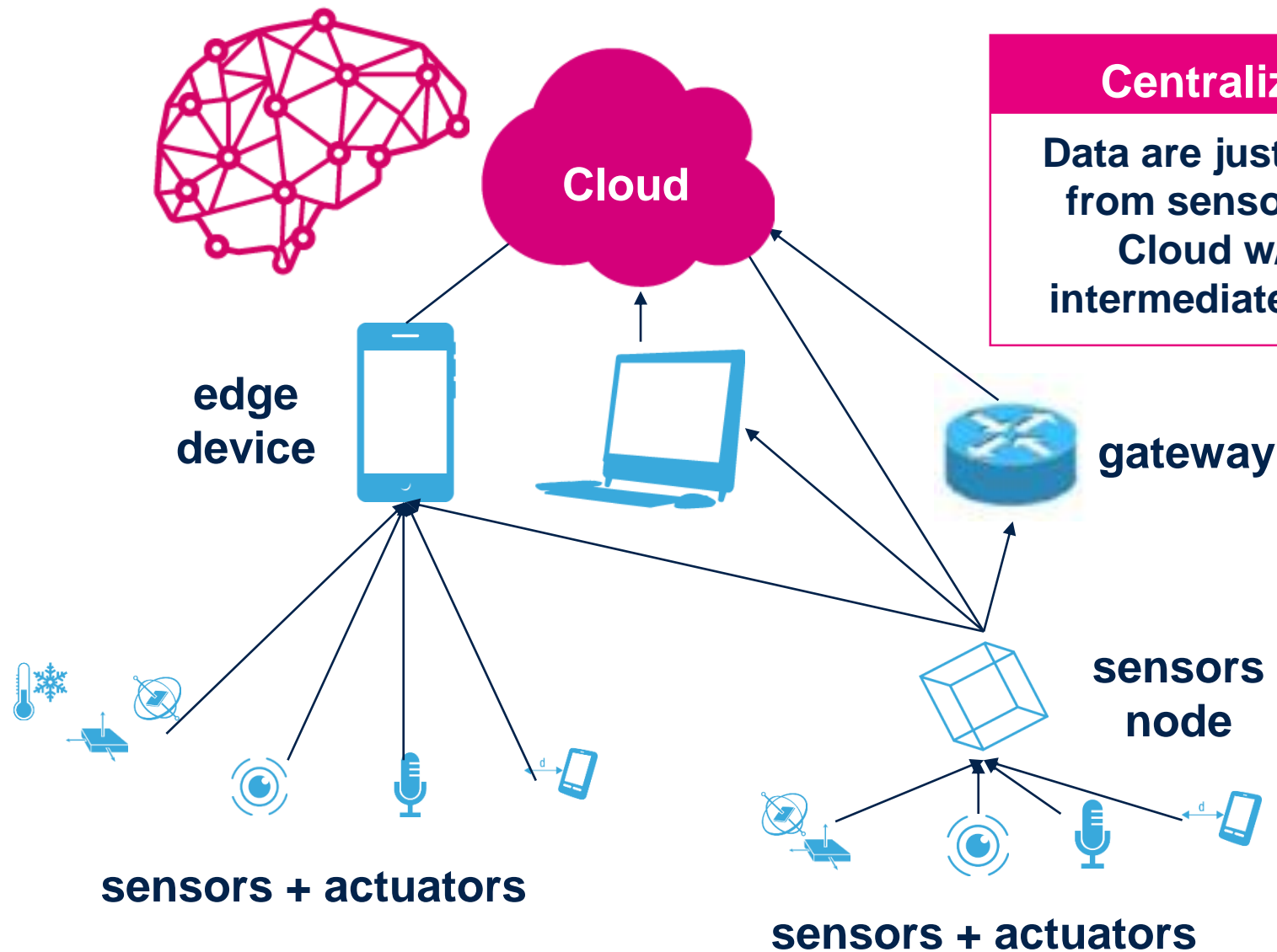







Edge AI is driven by IoT

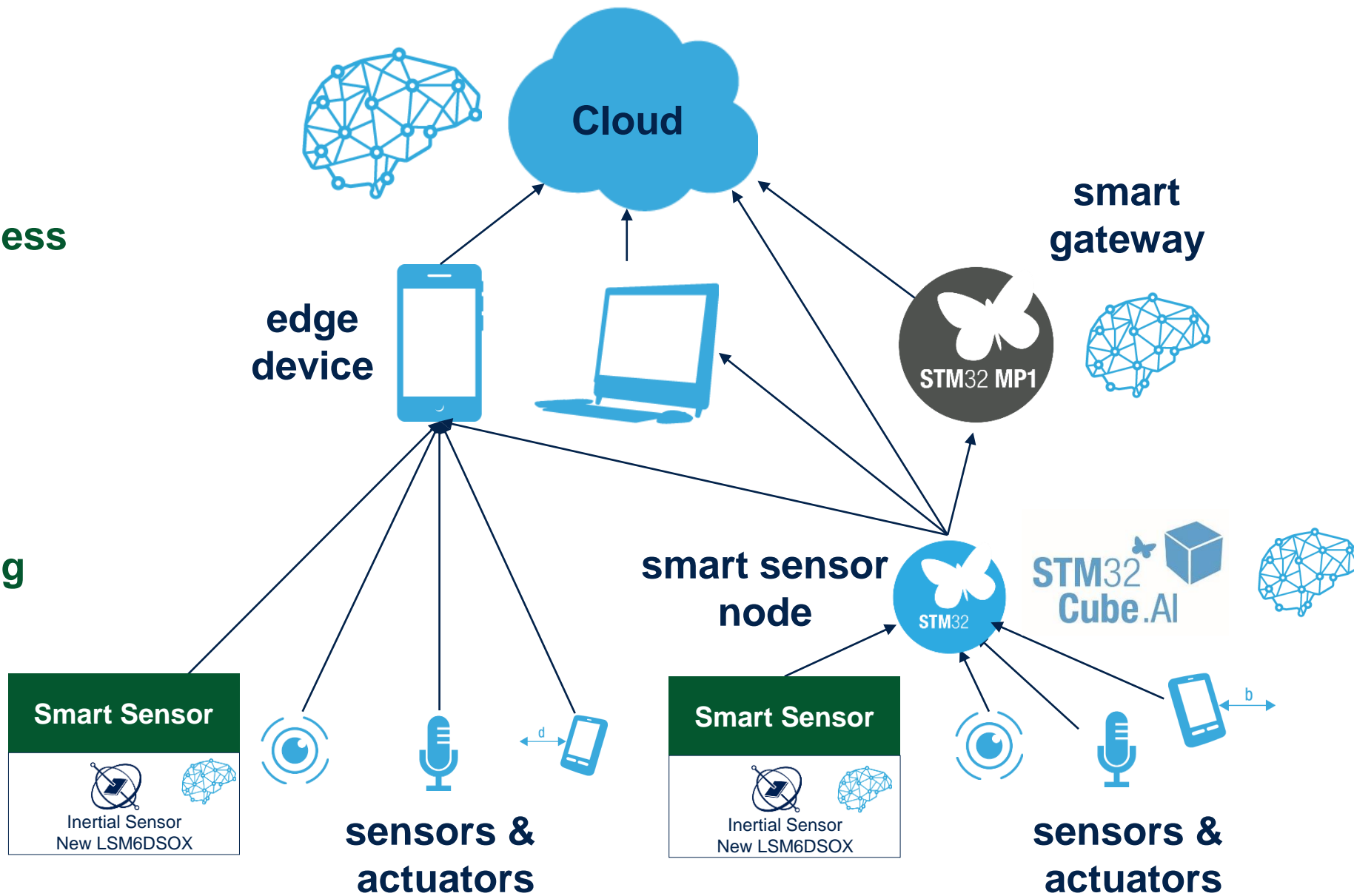


Centralized approach AI in the Cloud

- ↓ Responsiveness
- ↑ Bandwidth
- ↓ Privacy
- ↓ Security
- ↓ Energy Saving



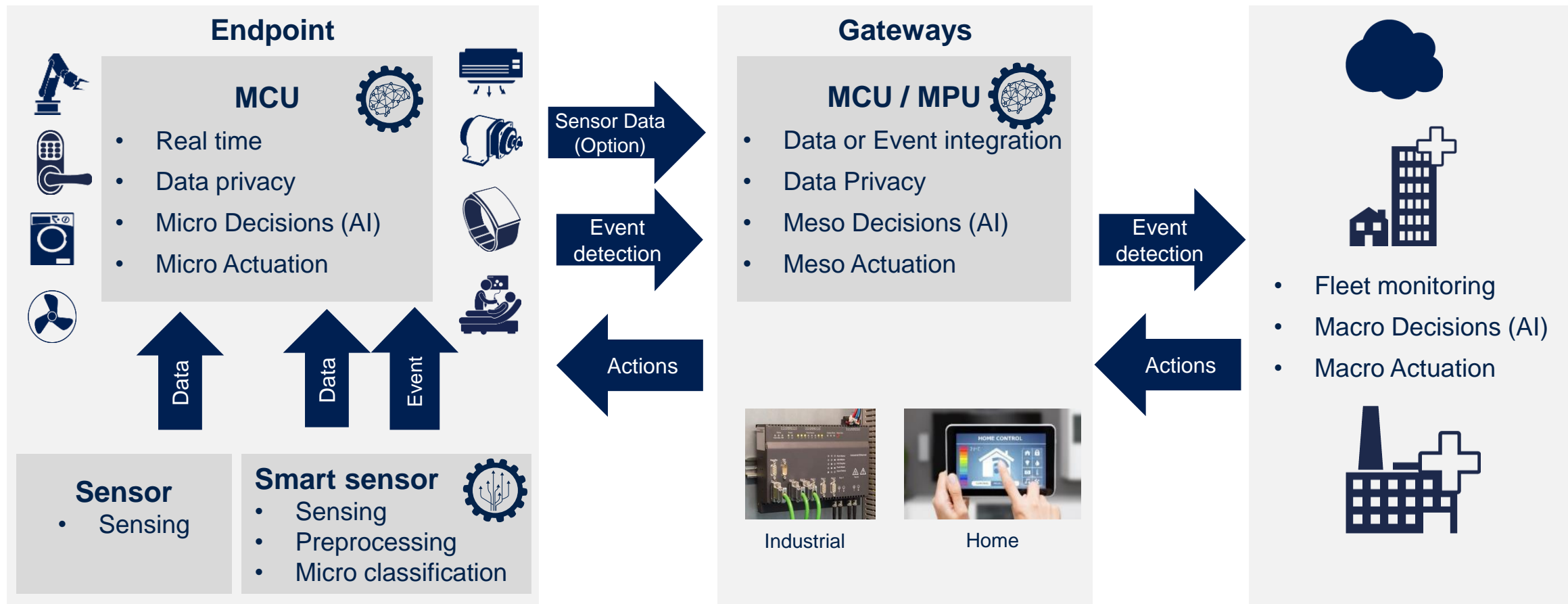
-  Responsiveness
-  Bandwidth
-  Privacy
-  Security
-  Energy Saving



Distributed AI from Edge to Cloud

Edge

Cloud



Machine level
decisions

Room or building-level
decisions

Cities or factory-level
decisions

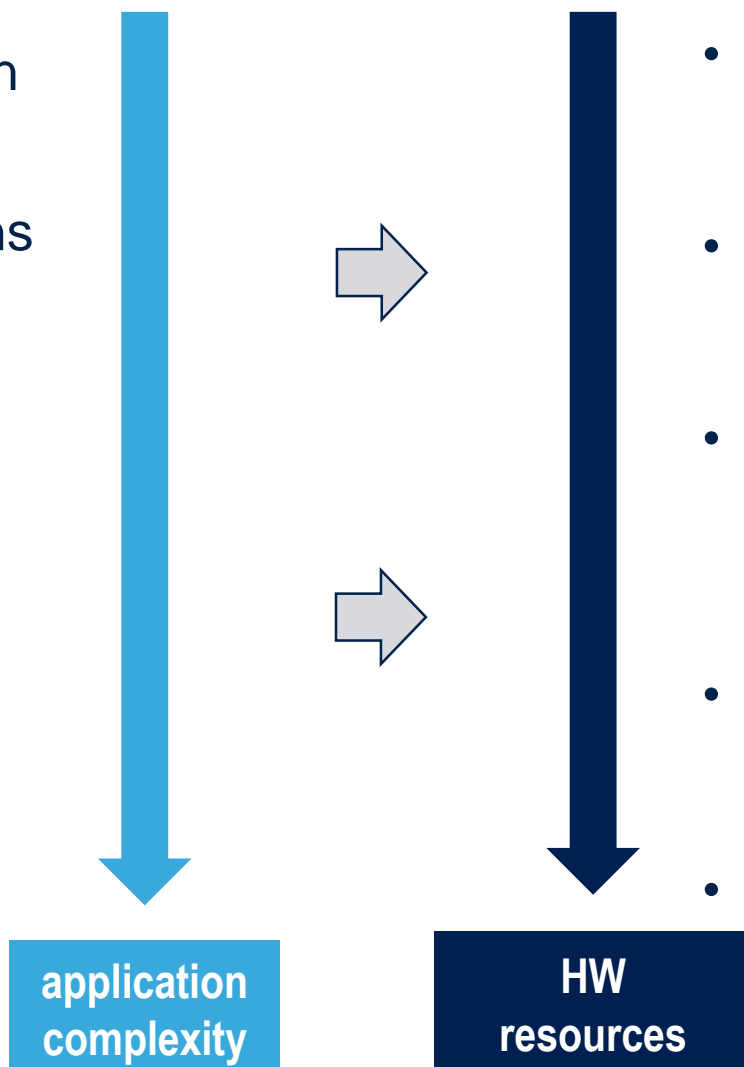
Edge AI computing technologies



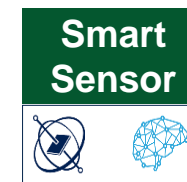
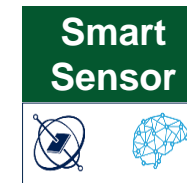
STM32 
Cube.AI

The right HW for the right application

- Simple event detection
- Multi-sensors decisions
- Activity recognition
- Context awareness
- Speech Recognition
- Computer Vision

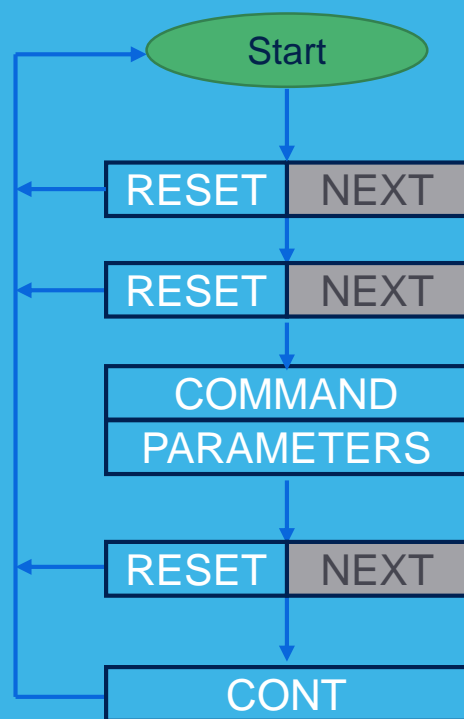


- Smart Sensors
- Smart MCU
- Smart Sensors & Smart MCU
- MCU + NPU
(Neural Processing Unit)
- Smart MPU

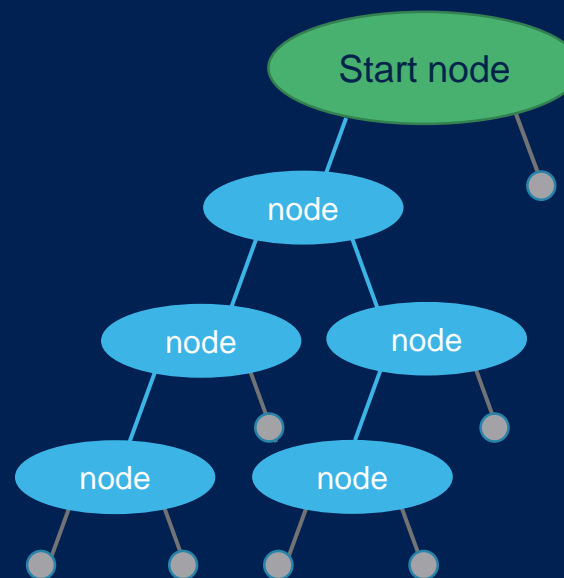


Add intelligence *INSIDE* the sensor

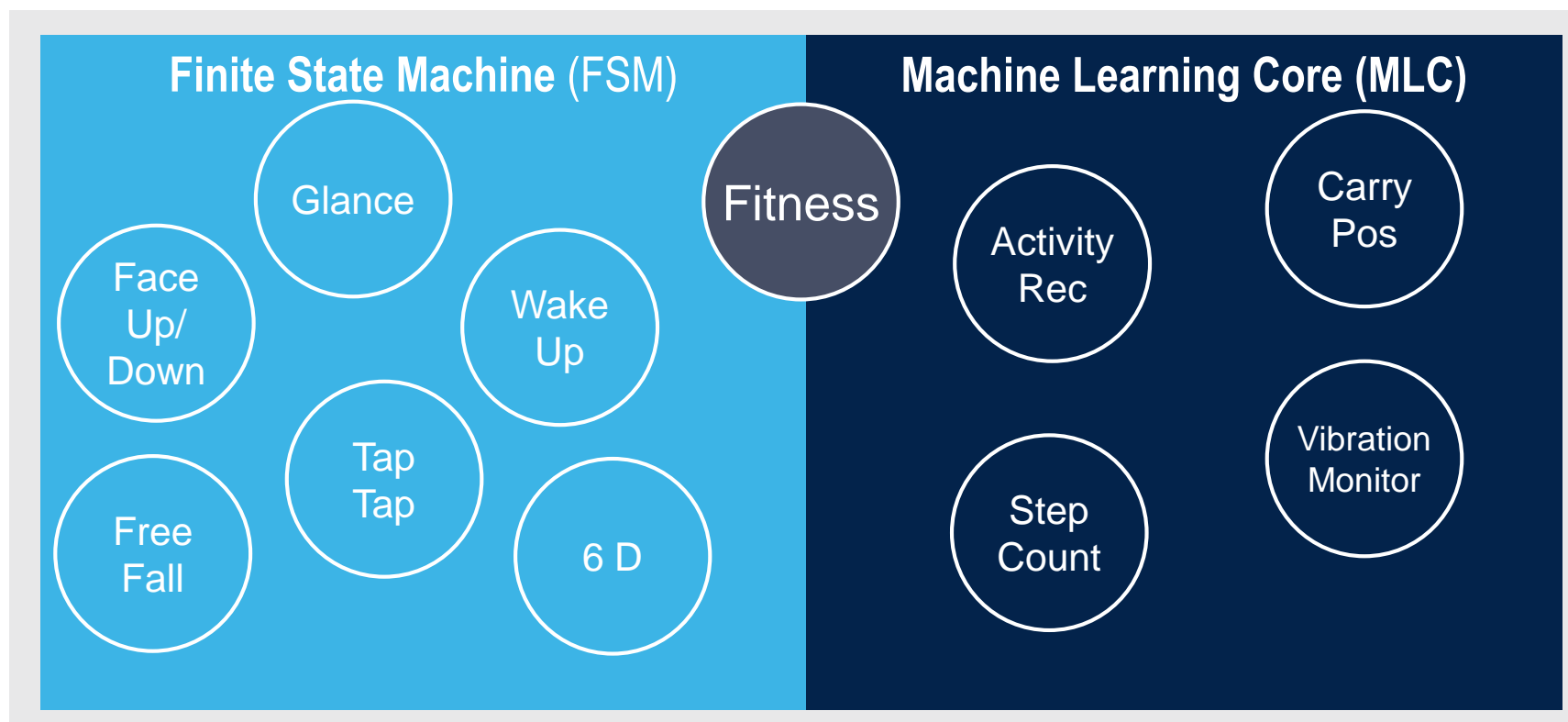
Finite State Machine

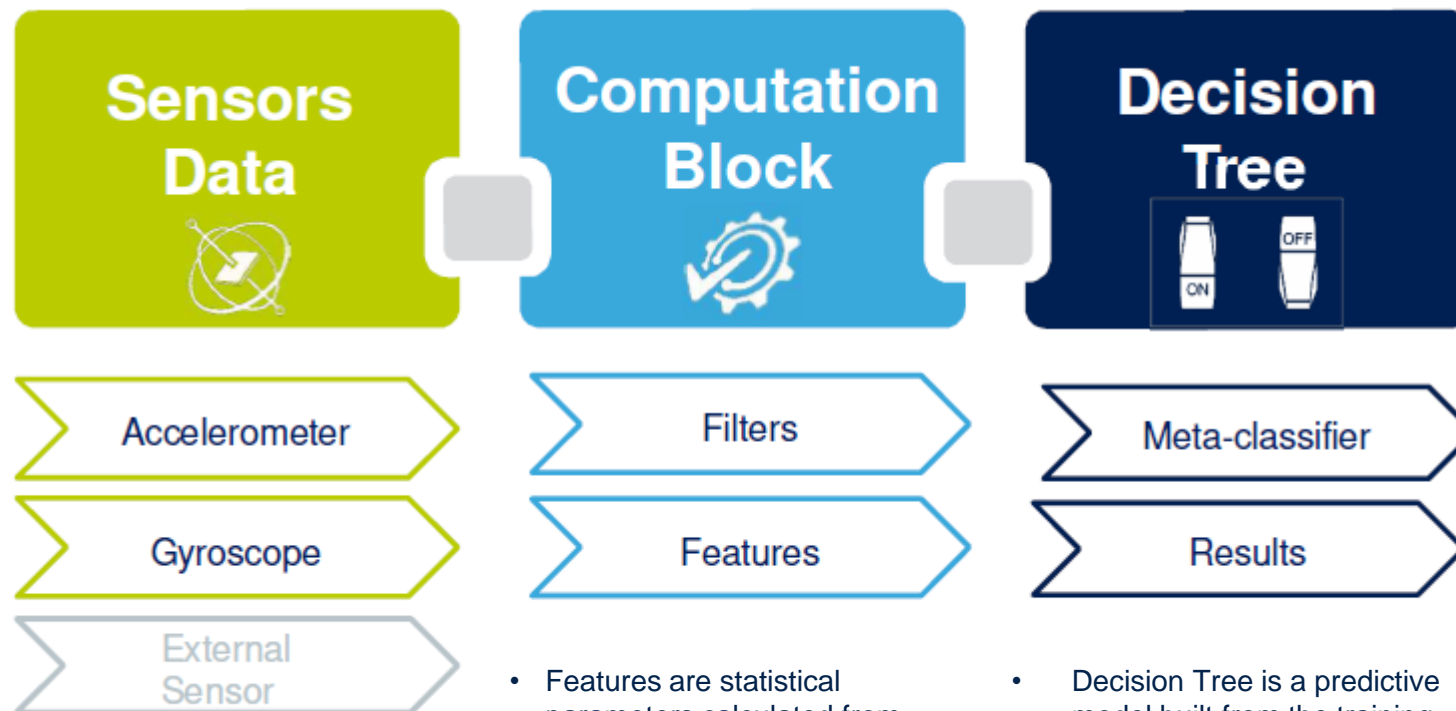


Machine Learning Core



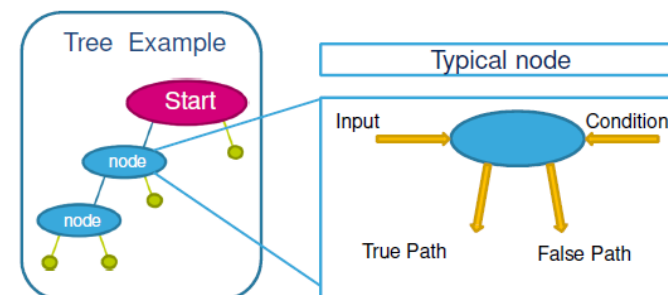
Finite State Machine and Machine Learning Core





- Features are statistical parameters calculated from
 - Input data (accX, accY, gyroX, gyroY, ...)
 - Filtered data (high-pass on accZ, band-pass on accY)
- Examples:
 - Mean
 - Variance
 - Energy
 - ...

- Decision Tree is a predictive model built from the training data



Activity recognition use case

10 to 1000 time energy saving by running MLC on Sensor vs. MCU/AP

How it works in 5 simple steps and with an intuitive use case:



User defines **Classes** to be recognized



Collect, clean and label data **Logs** according the classes



Define **Features** that best characterize the identified classes



Machine Learning tools generate program based on **Logs and Features**



Configure the LSM6DSOX and **run** the application



Capture data



Label data



Build decision tree



Embed decision tree

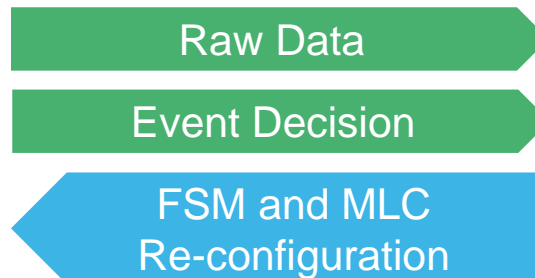


Process new data



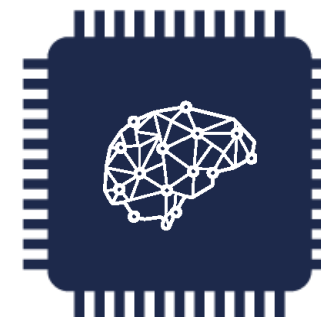
Smart Sensor + Smart MCU

Smart Sensor with Machine Learning Core



- Best ultra-low-power sensing at high performance:
 - 550 μ A (gyroscope and accelerometer)
→ 200 μ A less than closest competitor
 - 20~40 μ A (Accelerometer only for HAR)
- Efficient Finite State Machines: 3 μ A
- Configurable Machine Learning Core: 1~15 μ A

Smart STM32 second level of AI processing



Deep Learning
Neural Networks
Machine Learning

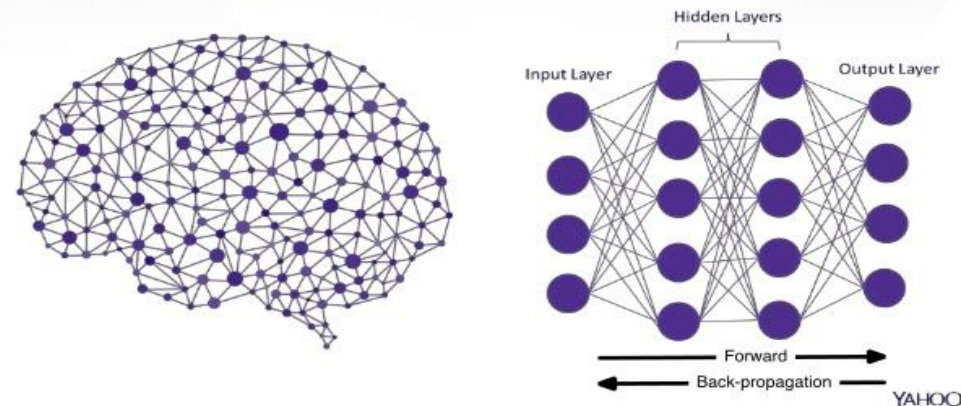
- More advanced and complex NNs
- Decisions on multiple sensors
- NN input can be sensor data and/or sensor Machine Learning decisions
- Multiple Neural Networks support
- Actuation & communication



Deep Learning (DL)

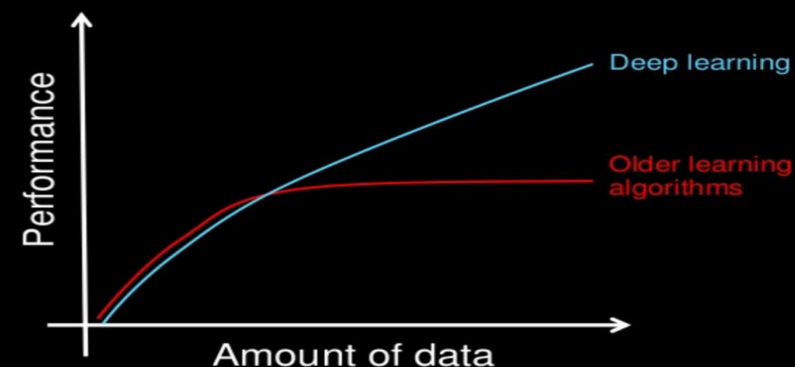
- Deep Learning is ML using Neural Networks
 - Inspired by biological neural networks
 - Deep because of having many intermediate learning steps
 - Lots and lots of data is required

Deep Learning



Advantages	Disadvantages
Autonomous Learning of data patterns & relationships	Large Datasets
High Accuracy	High Computational Requirements
Easy Improvement & Fine Tuning	Weak theoretical explanation
Adaptive Solutions	Black box (for most people)

Why deep learning



How do data science techniques scale with amount of data?

Why Deep Learning is so important

- Convolutional Deep Neural Networks outperform previous 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%

May 2015

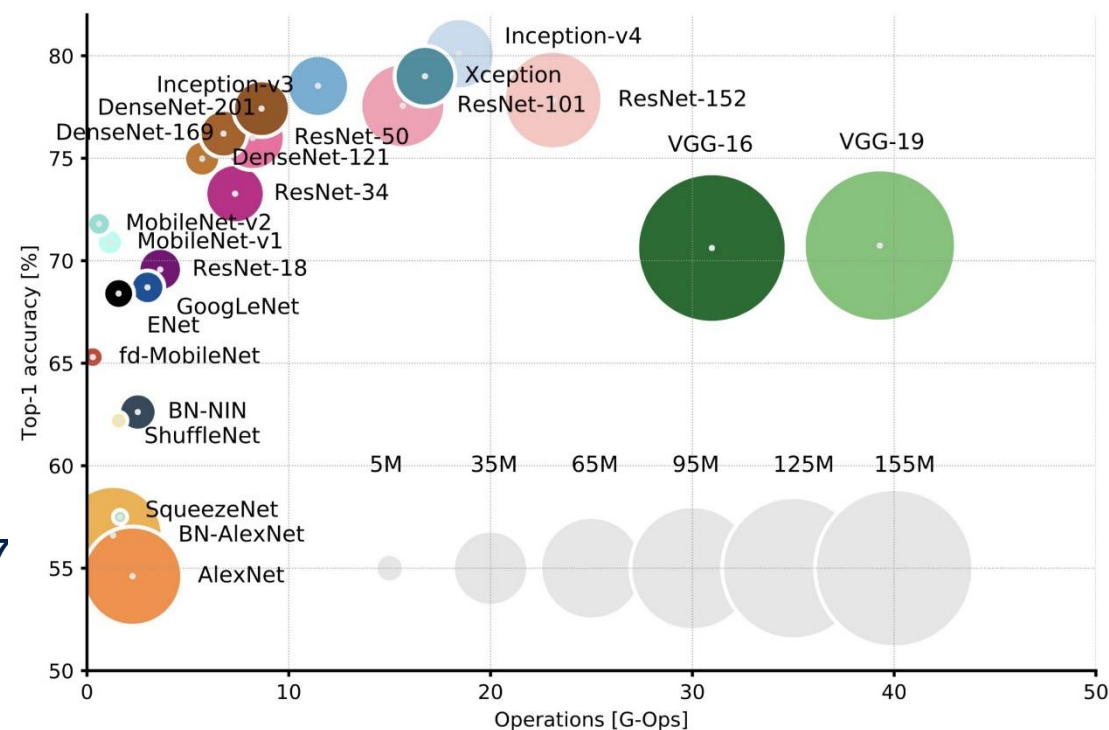
Neural Networks on MCU

A real challenge

- A MCU usually has:
 - Limited Non Volatile Memory (< 2 MB) → take care of # weights
 - Limited SRAM (< 1 MB) → take care of data size & activations
 - Low frequency (< 500 MHz): take care of # ops for time constrained applications
- The challenging task is to fit these constraints with good quality results !
- An example: image classification task
 - Gops should be very low
 - #weights should be very low
 - MobileNet-v2 is still high in terms of #weights and cannot run in real time, i.e., @30fps
 - We implemented FD-MobileNet (on 18 foods) on STM32H7
 - Memory footprint: 205 KB SRAM, 191 KB Flash !



Inference time 150 ms



The key steps behind Neural Networks



Neural Network (NN) Model Creation



Operating Mode

Capture data



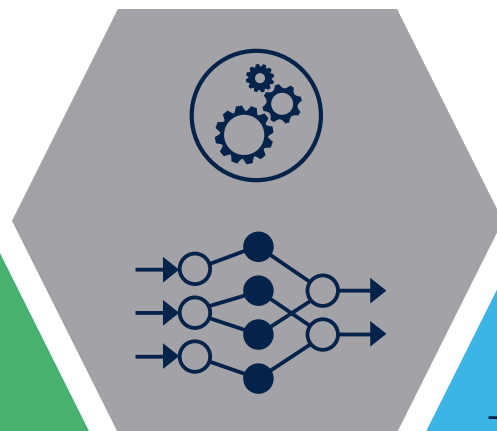
1

Train NN Model

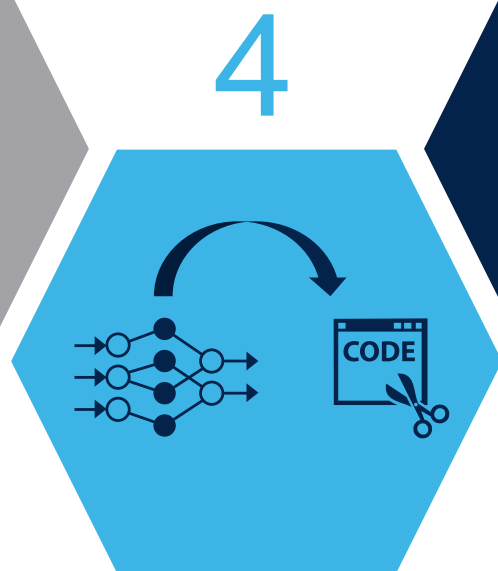


2

Clean, label data
Build NN topology



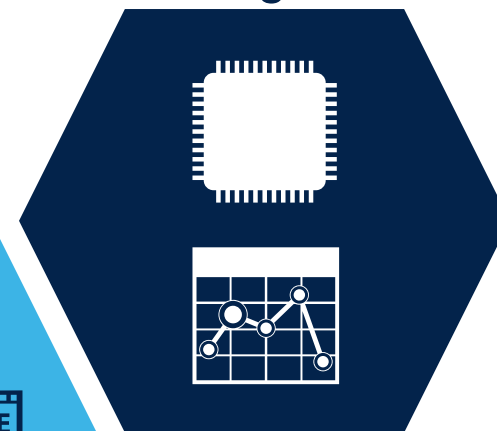
3



4

Convert NN into
optimized code for MCU

Process & analyze new
data using trained NN



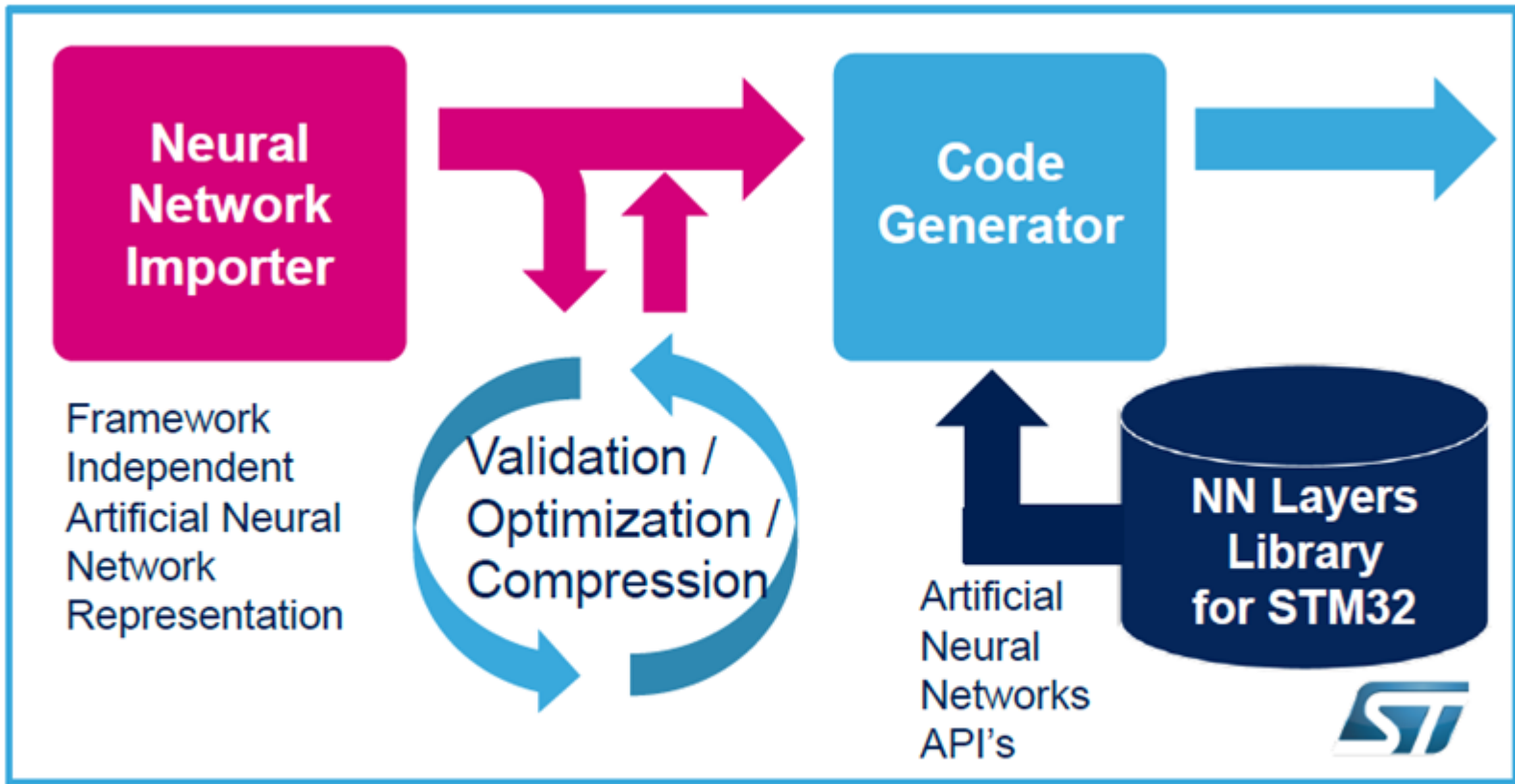
5

Neural Networks on MCU

The ST's way

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

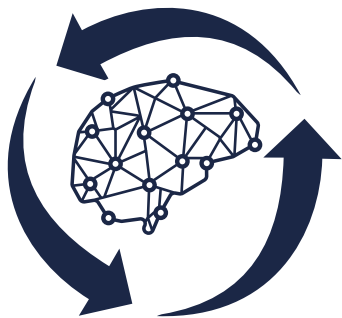
Optimize after initial deployment









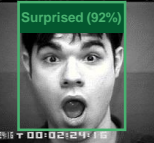
Continuous learning & Over-the-air update

STM32 
Cube.AI

Refine Neural Network to local conditions



- Learning at the edge to specialize NN to local sensor input: personalize to a particular user, home/factory environment...
- Supervised learning: user inputs feedback to re-enforce classification or regression output
- Unsupervised learning: classification with high probability is used for re-enforcement, output labels are guessed from inferences in same temporal window

	Activity monitoring
	Home Audio Event Detection
	Food recognition
	Touch gesture recognition
	Motor anomaly detection
	Written character recognition
	Face expression

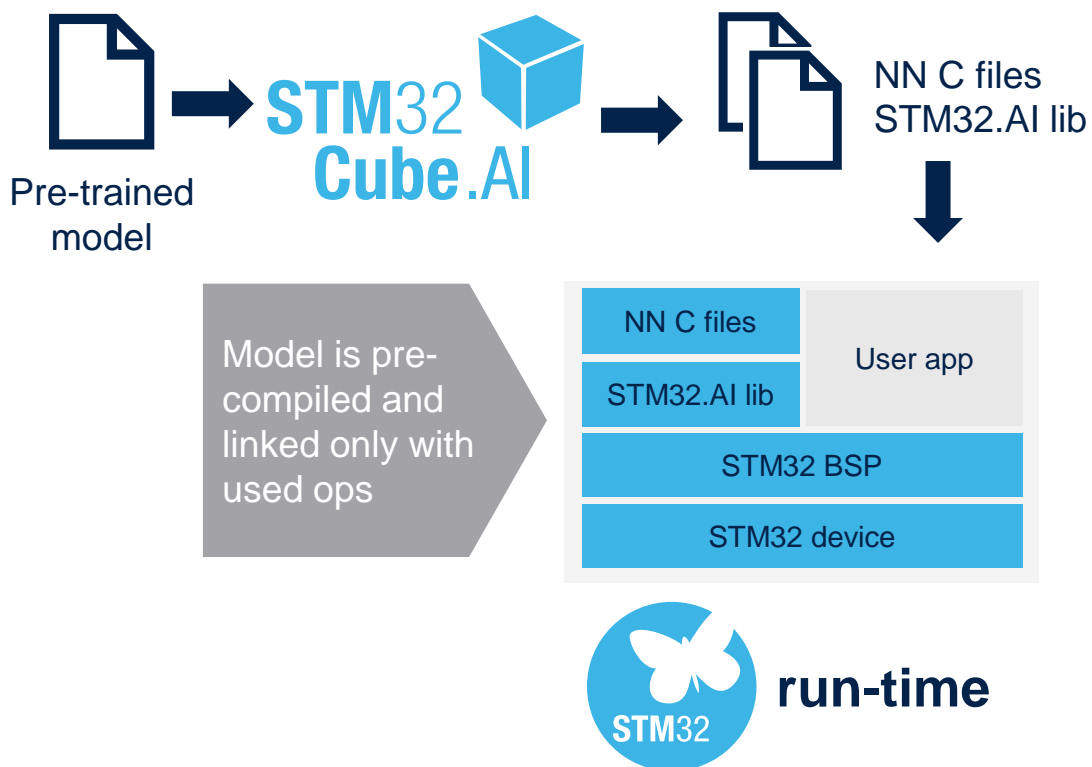
Refine Neural Network from learning of all local conditions



- Learn at the edge from different conditions, send learning cumulated over batch of data to the Cloud regularly
- Cloud cumulates learnings from all sensors: builds learning from diverse dataset without full raw data transfer over connectivity
- Generalization is ensured by federated learning of many diverse examples

	Activity monitoring
	Home Audio Event Detection
	Food recognition
	Touch gesture recognition
	Motor anomaly detection
	Written character recognition
	Face expression

NN update over the air



Update weights

Using STM32Cube.AI v5.0.0, model weights are stored in specific table of NN C files

- Fine tune weights (when learnt in the Cloud), independently of the FW

Update topology

From STM32Cube.AI v6.0.0, Neural Networks files and AI library are stored in specific sections. Topology can be updated without full FW update

- Add a new class
- Add an extra layer
- Add new operators

Application markets

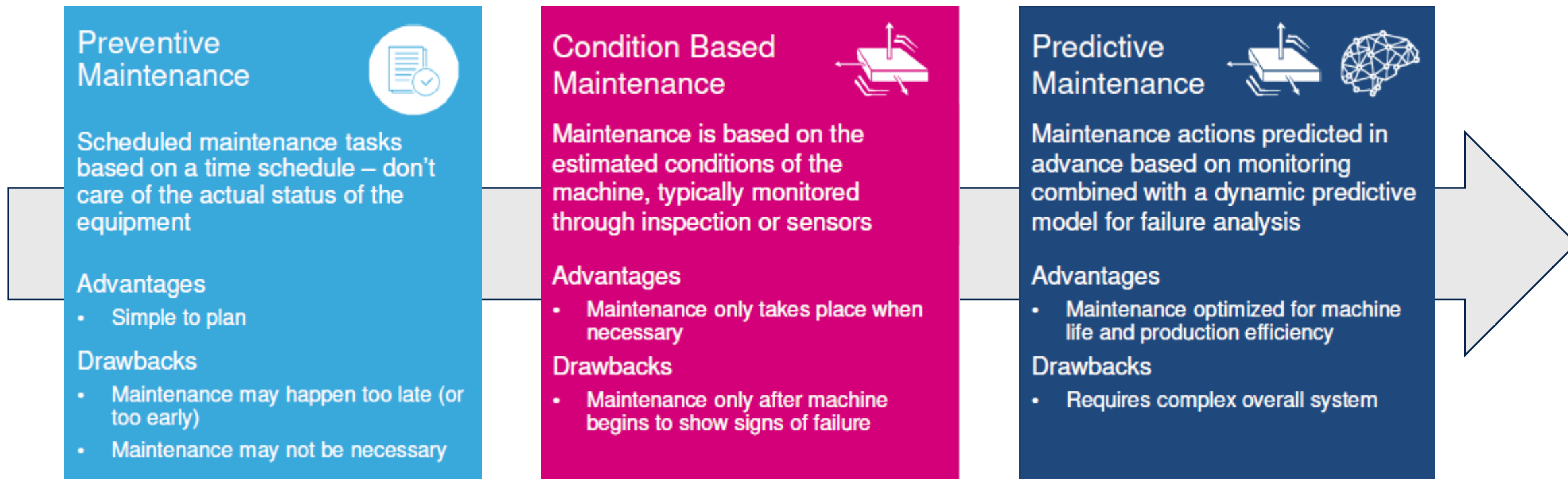


Condition monitoring and predictive maintenance

STM32 
Cube.AI

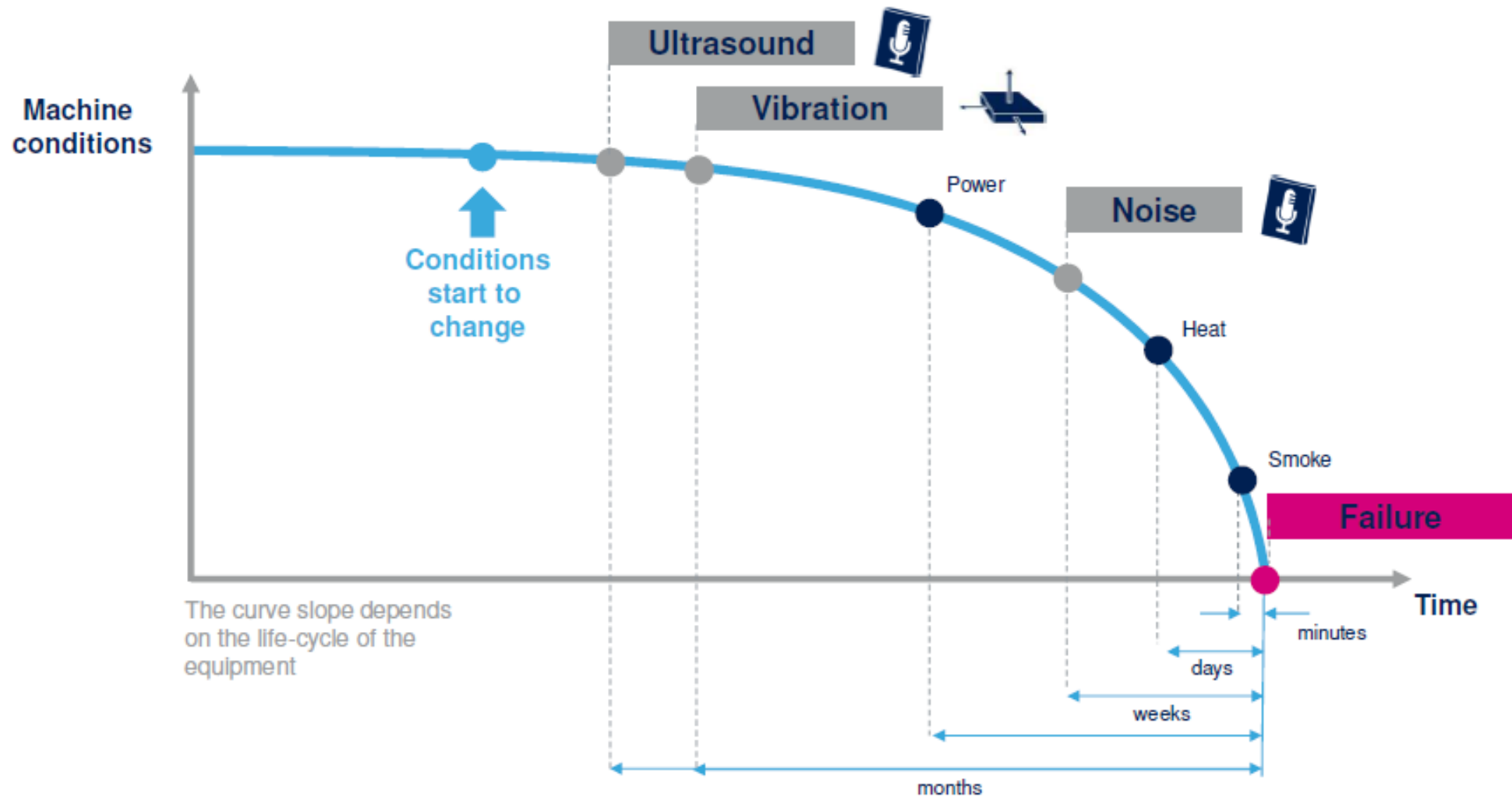
Predictive Maintenance

A Smart Industry hot topic



Predictive Maintenance

Microphones and inertial sensors



Predictive Maintenance

Benefits



Reduced lost production time

Maintenance on the production line only when needed and at the optimal time

Longer machine lifetime/lower effective cost

Replacing the minimum amount of parts before failure causes damage to others

Faster and more efficient repair

Optimized workers interventions and minimum labor for parts replacement

Increased safety

Prevents failures that could be dangerous for workers before they happen



*Aggregated figures from different sources
(Accenture, McKinsey, ST)*

Condition monitoring & predictive maintenance applications

Industrial

- ❖ Manufacturing and Process Automation
- ❖ Power and energy
- ❖ Home appliances and Smart Building Automation

Automotive and Transportation

- ❖ Vehicles
- ❖ Railways
- ❖ Infrastructures

Home Appliance

- ❖ Lighting
- ❖ Washing Machine
- ❖ Vacuum cleaner, Air conditioning

New equipment
(greenfield)

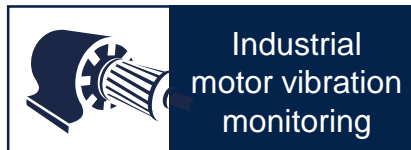


Integration possible
with power supply
and existing sensors

In-field maintenance
(retrofit)



Battery-powered
simplifies
installation



Industrial
motor vibration
monitoring



Engine TTF
vibration
monitoring



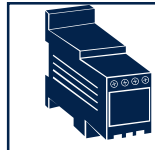
Motor current
monitoring



Crankshaft
rotation
monitoring



Pipe flow
monitoring



Temperature,
humidity, gas
monitoring



acoustic
monitoring

Application markets



Sound event and context awareness

STM32 
Cube.AI

Sound-based applications

Sound event classification

- ❖ Home alarms
- ❖ Glass break
- ❖ Machine anomaly

Sound event detection

- ❖ Vehicle count based on Doppler
- ❖ Sound direction detection with multiple MICs

Sound wake-up

- ❖ Sound detection to wake up more advanced algorithms

Context awareness

- ❖ Indoor/outdoor environment
- ❖ Vehicular environment
- ❖ Factory activity
- ❖ Human presence detection



Complex audio front-end not always required for sound detection

	Vehicle counting
	Intrusion detection
	Engine break detection
	Audio scene classification
	Building occupancy
	Crowd panic detection
	Wake-up / orient camera based on sound

Application markets






Computer Vision

STM32 
Cube.AI

Computer Vision




Processing requirements

Low




- Static low-res images
- Known object position
- Good light conditions

Medium

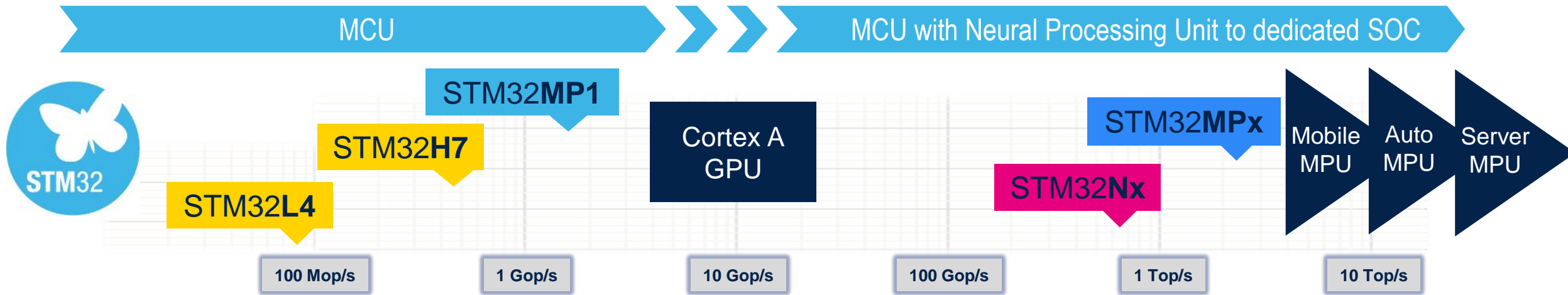




- Low frame rate / QVGA images
- Open environment
- Medium light conditions

High

- High rate Video
- High resolution
- Adaptive light conditions



Vision based applications

Computer Vision for voice interfaces

- ❖ Gaze detection
- ❖ Gesture recognition
- ❖ Person direction detection

Predictive maintenance

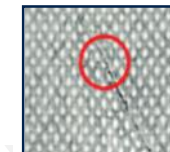
- ❖ Smoke detection
- ❖ Product defect detection
- ❖ Spilled liquid detection

Visual wake-up

- ❖ Person detection to wake up more advanced algorithms

Classification and recognition

- ❖ Recognize pests, weeds, disease in fields
- ❖ Characters and digits recognition
- ❖ Ingested camera imaging
- ❖ Simple image classification
- ❖ Texture, fabric recognition
- ❖ Visual biometrics



Defect detection



Crop disease detection & classification



Thermal camera



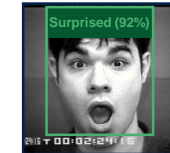
Meter aftermarket



Gesture recognition



Person detection



Face expression

IoT object often rely on batteries, as main power is mostly not available



Replacing batteries is not sustainable in IoT, low power is a MUST

Servitization



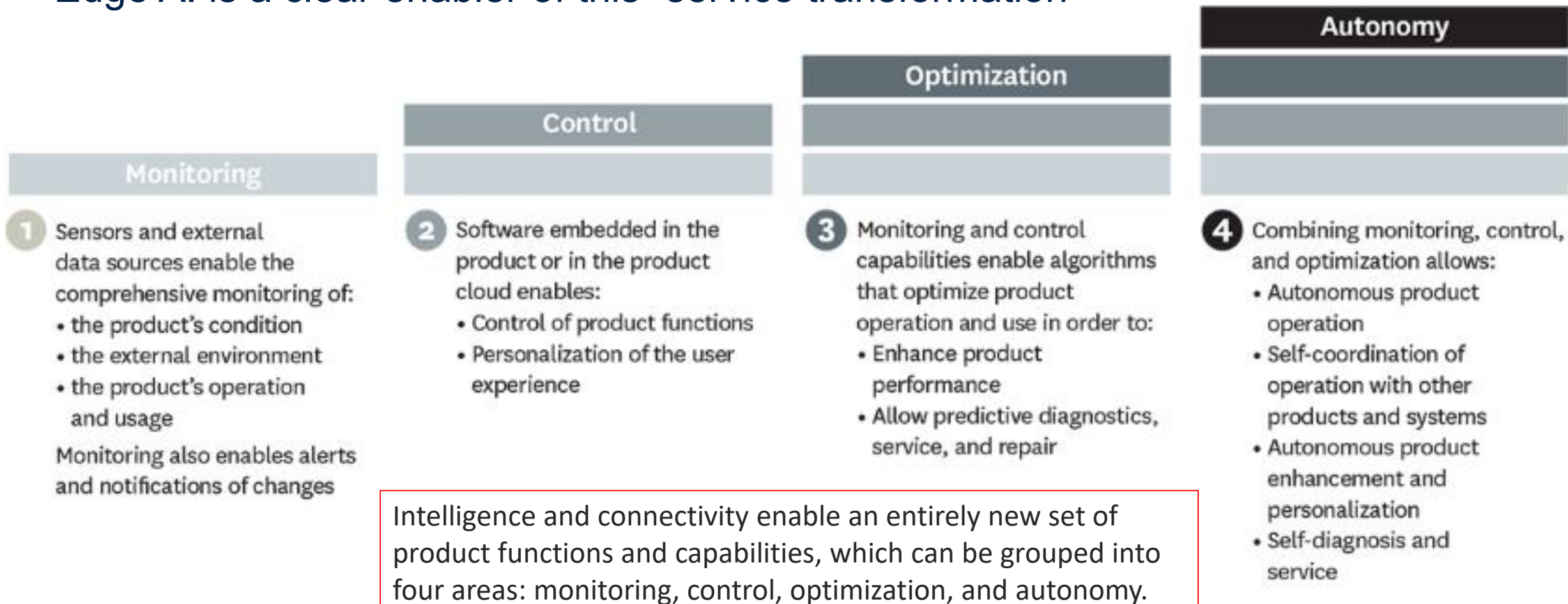
STM32 
Cube.AI

- **Servitization** means the joint offering of products and services, or better “product-service systems”
- Companies can better create value by moving from the sale of products to the sale of systems consisting of products and services
- **Servitization** enables companies to differentiate their offer from that of their competitors and at the same time increase customer loyalty over time
- Some examples:
 - IBM selling Cloud services (like *Watson* and *Quantum Computing*), instead of selling hardware (mainframes)
 - Rolls-Royce selling “power-by-the-hour“, instead of selling aircraft engines



Edge AI and servitization

- Edge AI is a clear enabler of this “service transformation”





FEDERAZIONE NAZIONALE
IMPRESSE ELETTRONICHE
ED ELETTRONICHE



Thank you

© STMicroelectronics - All rights reserved.

The STMicroelectronics corporate logo is a registered trademark of the STMicroelectronics group of companies. All other names are the property of their respective owners.



life.augmented