# tinyML® Summit

*Miniature dreams can come true…*

**March 28-30, 2022 | San Francisco Bay Area**

www.tinyML.org

# Brains into sensors with AI in the Edge

Andrea Onetti
Executive Vice President MEMS Sensors Sub-Group
Analog, MEMS and Sensors Group

STMicroelectronics

# The MEMS journey



| Offline era | Online era | Onlife era |
|---|---|---|
| **2000** | **2010** | **2020** |
| **A paradigm change in the man-machine interface** | **Sensor's proliferation and connections to Cloud** | **The fusion of technology and life** |
| MEMS technology: from a concept to a product | Performance improvement and technology fusion | Standalone devices able to sense, process and take action |

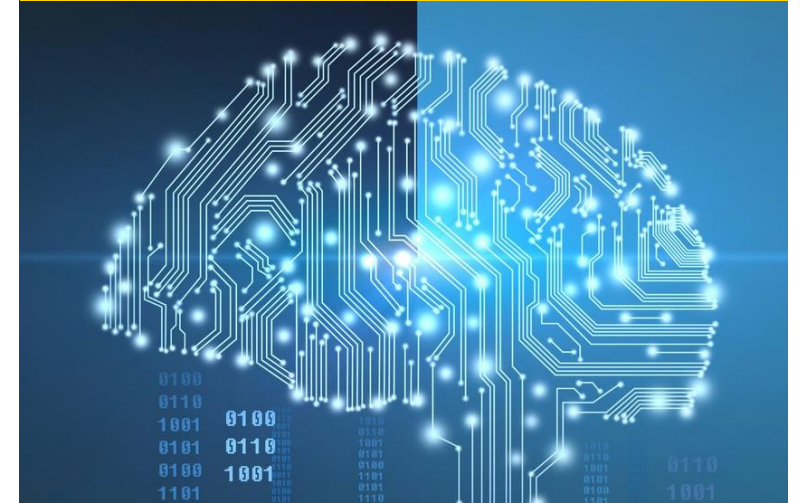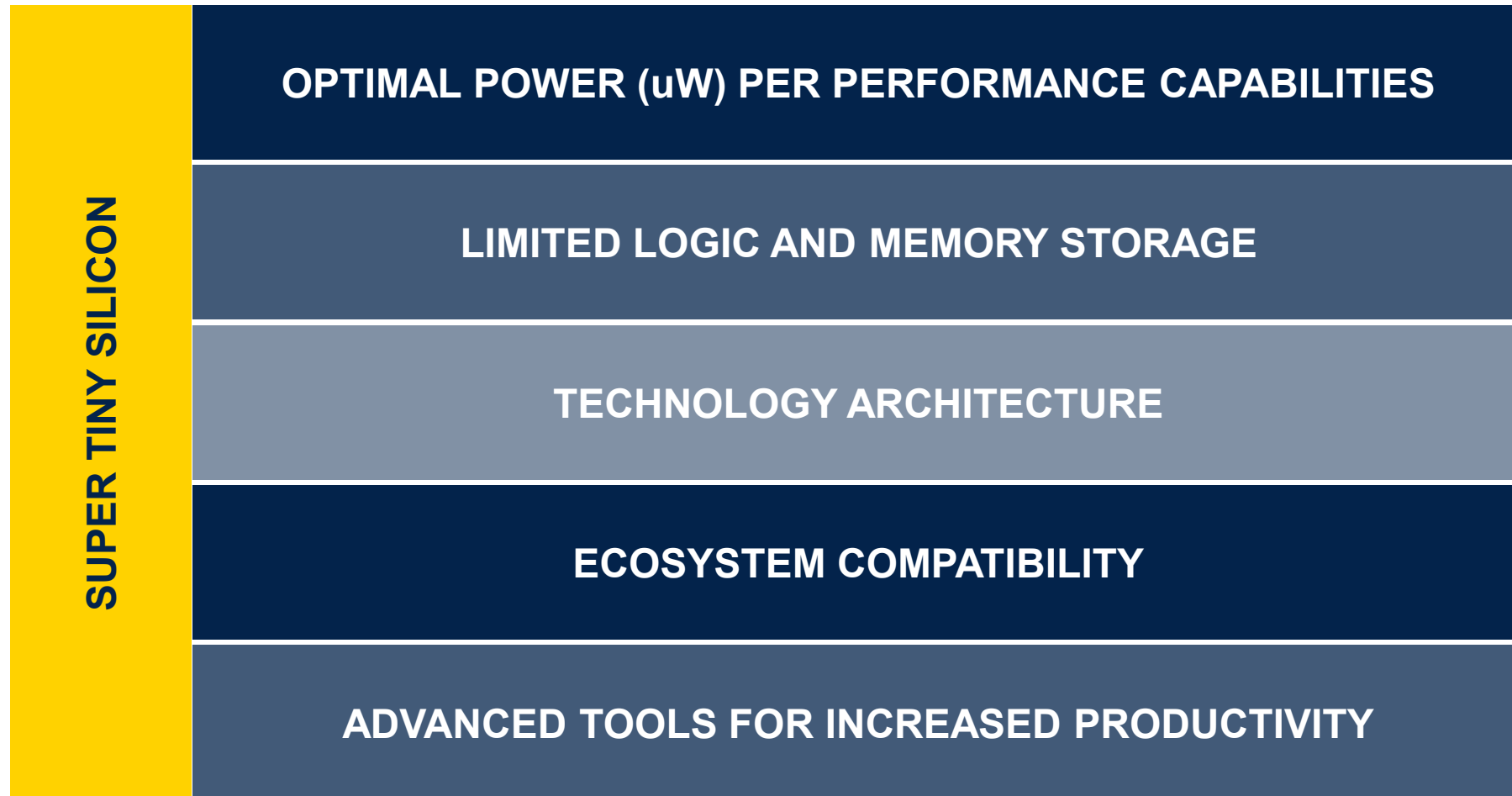| Offline era | Online era | Onlife era |
|---|---|---|
| **Fragmented** | **Connected** | **Trained** |
|  |  |  |
| The simplest configuration: independent systems | Intertwined nodes enable efficient data exchange | Edge AI local decision making with maximum privacy |

# Industry 5.0 challenges

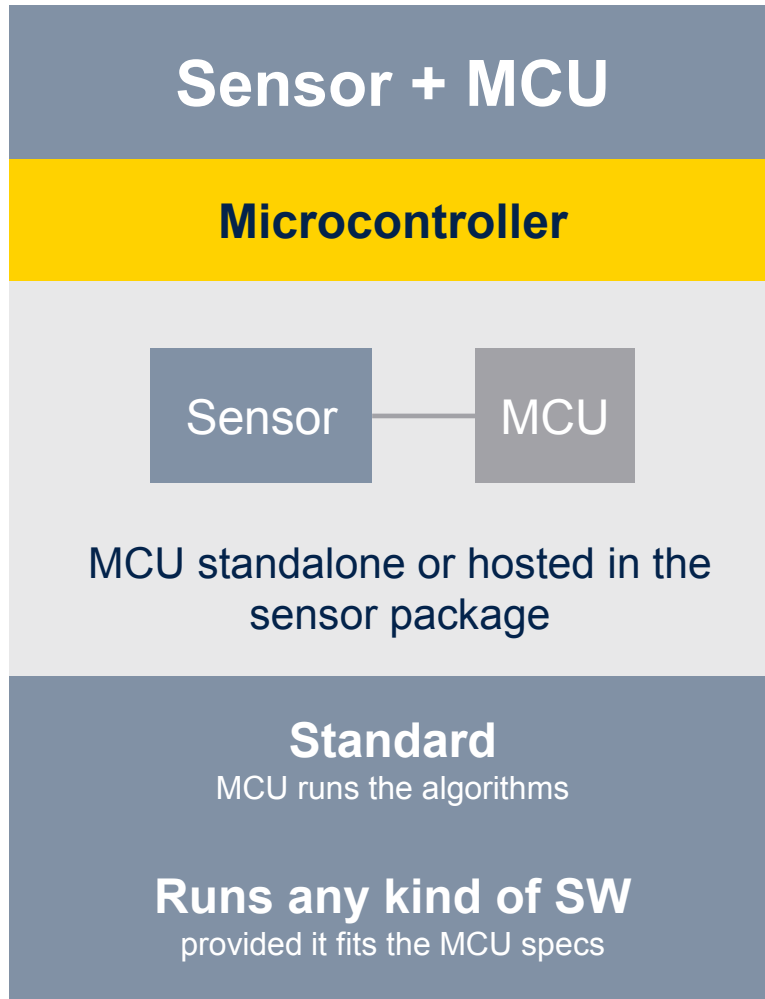# Sensor's semiconductors challenges for Edge AI

**SUPER TINY SILICON**

- OPTIMAL POWER (uW) PER PERFORMANCE CAPABILITIES
- LIMITED LOGIC AND MEMORY STORAGE
- TECHNOLOGY ARCHITECTURE
- ECOSYSTEM COMPATIBILITY
- ADVANCED TOOLS FOR INCREASED PRODUCTIVITY

life.augmented

# Migrating intelligent processing
## From "on the Edge" to "in the Edge"

| Sensor + MCU | rPU | ISPU |
|---|---|---|
| **Microcontroller** | **reconfigurable Processing Unit** | **Intelligent Sensor Processing Unit** |
| Sensor — MCU | rPU+ Sensor — MCU | ISPU + Sensor — MCU |
| MCU standalone or hosted in the sensor package | rPU integrated in the sensor ASIC | ISPU integrated in the sensor ASIC |
| **Standard** <br> MCU runs the algorithms <br><br> **Runs any kind of SW** <br> provided it fits the MCU specs | **Optimized** <br> reconfigured through register setting <br><br> **Constrained** <br> runs same model/mapping | **Programmable** <br> dedicated instruction set <br><br> **Runs several AI algorithms** <br> Full precision to 1-bit NN |

## DSP for real-time processing and Artificial Intelligence
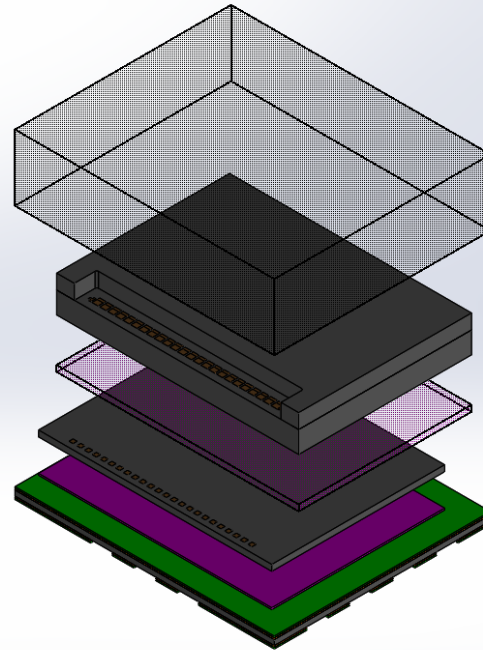
**Small area**
down to 8 kgates

**Standard package**
3 x 2.5 x 0.83 mm

**RAM based**
40 kiloBytes (program + execution)



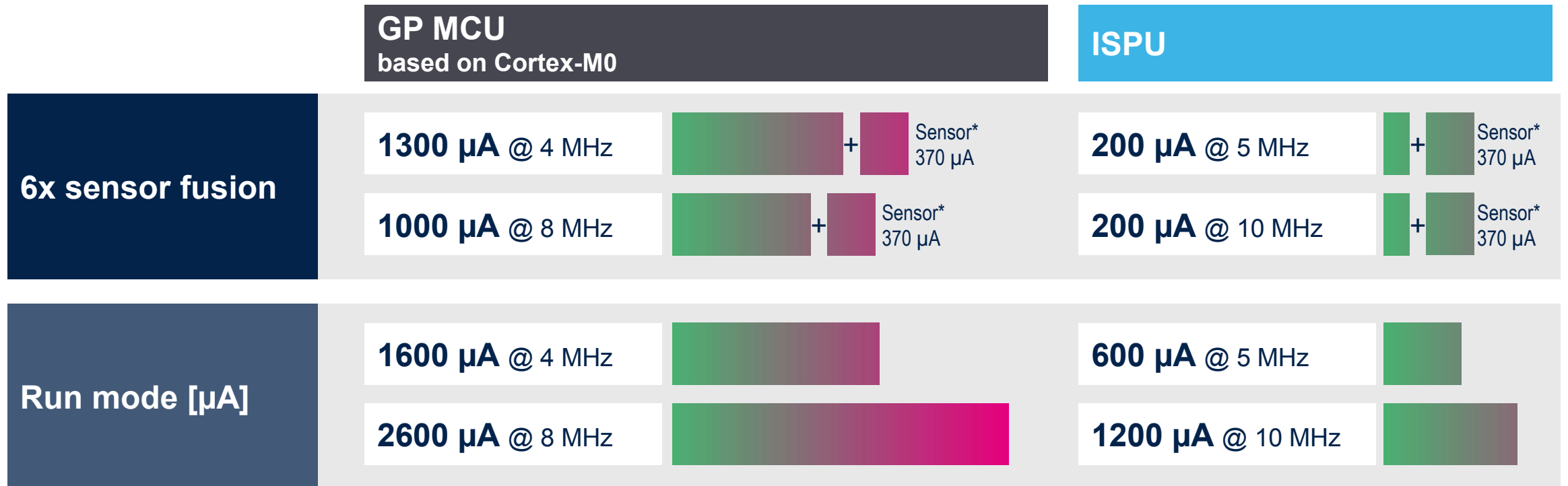**Full precision**
Floating Point Unit

**Binary Neural Network**
convolution acceleration

**Optimization**
Power consumption vs. performance

life.augmented

# Optimization: power consumption vs performance

**5x less current consumption for sensor fusion on ISPU than on M0**
**Below 600 µA for sensor fusion in the edge**

| | GP MCU based on Cortex-M0 | | ISPU | |
|---|---|---|---|---|
| **6x sensor fusion** | **1300 µA** @ 4 MHz | + Sensor* 370 µA | **200 µA** @ 5 MHz | + Sensor* 370 µA |
| | **1000 µA** @ 8 MHz | + Sensor* 370 µA | **200 µA** @ 10 MHz | + Sensor* 370 µA |
| **Run mode [µA]** | **1600 µA** @ 4 MHz | | **600 µA** @ 5 MHz | |
| | **2600 µA** @ 8 MHz | | **1200 µA** @ 10 MHz | |

*Accelerometer + Gyroscope low-power mode @ ODR 104 Hz

# Binary Neural Network (BNN) on ISPU

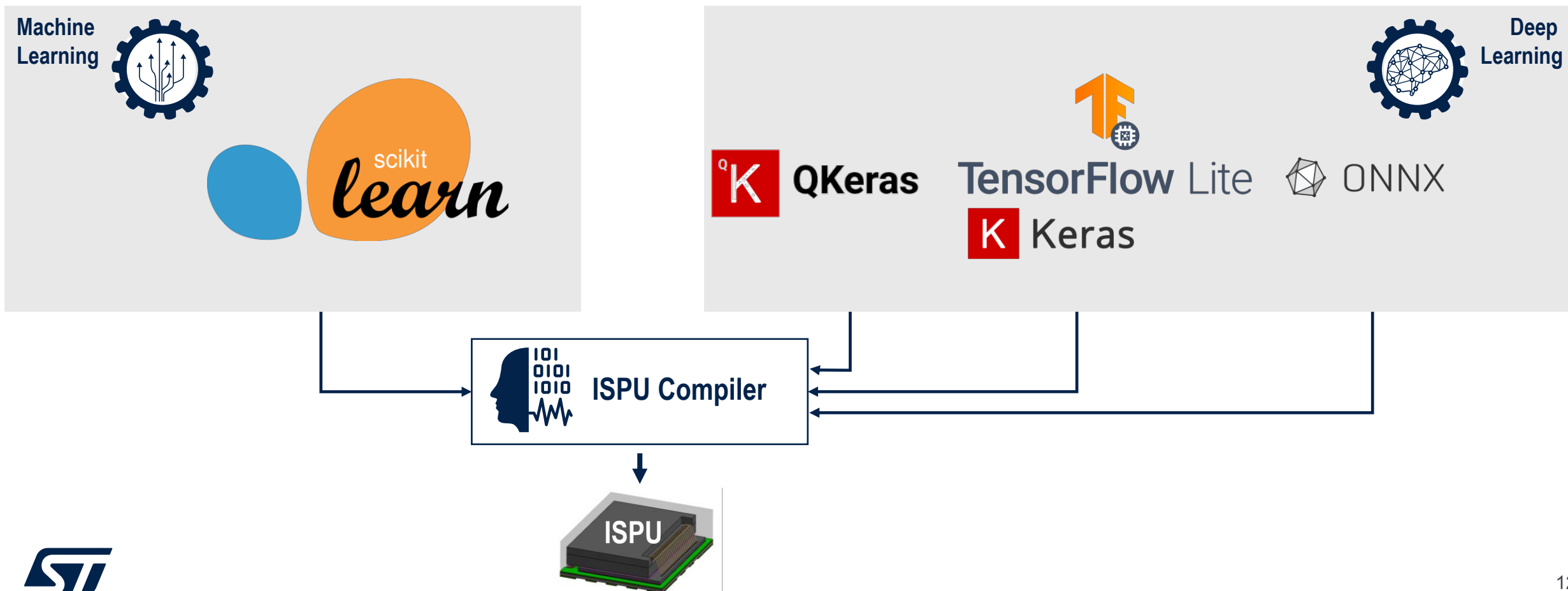**BNN on ISPU delivers over 10x better performance than floating point**

| Microbenchmark* [single dense layer 128x64] | Floating point | | Full BNN | | BNN Improvement |
|---|---|---|---|---|---|
| **Size**** | **106,324** Bytes | | **10,372** Bytes | | **10.3x** |
| **Execution time** | **74,210** Cycles | | **4,366** Cycles | | **17.0x** |

\* Kernel size = 128 and number of kernels = 64. ISPU set at 5MHz frequency
\*\* Full Application size: data + code + internal buffers + system libs

# Hybrid Binary Neural Network (BNN) on ISPU

## ISPU makes solutions ready for Onlife with faster and smaller algos

| Fan blade condition monitoring algorithm | Floating point | Hybrid* BNN QKeras | Hybrid* BNN Improvement |
|---|---|---|---|
| **Size**** | 107,200 Bytes | 11,404 Bytes | **9.4x** |
| **Execution time** | 246,170 Cycles | 194,380 Cycles | **1.3x** |

\* Some layers are floating point activations with binary weights, some are fully binarized (weights and activations). ISPU set at 5MHz frequency

\*\* Full Application size: data + code + internal buffers + system libs

# In-sensor Machine Learning & Deep Learning

**ST ISPU delivers more options and greater freedom**

| | |
|---|---|
| | **Very constrained silicon area for logic and RAM** <br><br> **No Flash memory** |
| | **Ultra-low power consumption (µW envelope)** |
| | **Easily programmable with AI commercial models** — NANOEDGE AI |
| | **Interoperates with** Keras, QKeras, TensorFlow Lite, ONNX, scikit learn |
| **?** | **And?** |

# The B.E.T. benchmark

**B**YTES
Amount of data transferred from sensor to cloud

**E**NERGY
Total system power consumption

**T**IME
From event to reaction: make local analysis cuts reaction time

# An example: the robotic arm handling

| Parameters | Offline | | Online | | Onlife | |
|---|---|---|---|---|---|---|
| | | **Benefits** | | **Benefits** | | **Benefits** |
| **Byte saving** (transferred from sensor to cloud) | | No data transfer | | | | No data to be stored or transferred |
| **Energy saving** (total consumption) | | | | Wafer waste reduced but data stored and processed on cloud | | AI/ML processing in the edge |
| **Time saving** (from event to reaction) | | | | Time to reaction reduced but still too long | | Machine Learning approach for failure detection / prediction |
| **OUTCOME** | 1 lot (25 wafers) wasted + machine calibration time | | 1 or 2 wafers wasted + machine calibration time | | No wafer wasted | |

# "In" the Edge: towards a new ecosystem

## What's need to be explored together for ISPU?

Ecosystem revision for tools, algorithms, and quantization procedures in sensors assets

Development of new benchmarks, and design tools to serve this innovation

Raise productivity and achieve synergies within the embedded developer community

# Takeaways

ISPU is real: global launch in 2022

ISPU is sustainable

ISPU empowers 10M+ C language developers in using AI in the Edge

ISPU makes Onlife possible

# Our technology starts with You

🌐 Find out more at www.st.com

life.augmented

Systems where sensors live

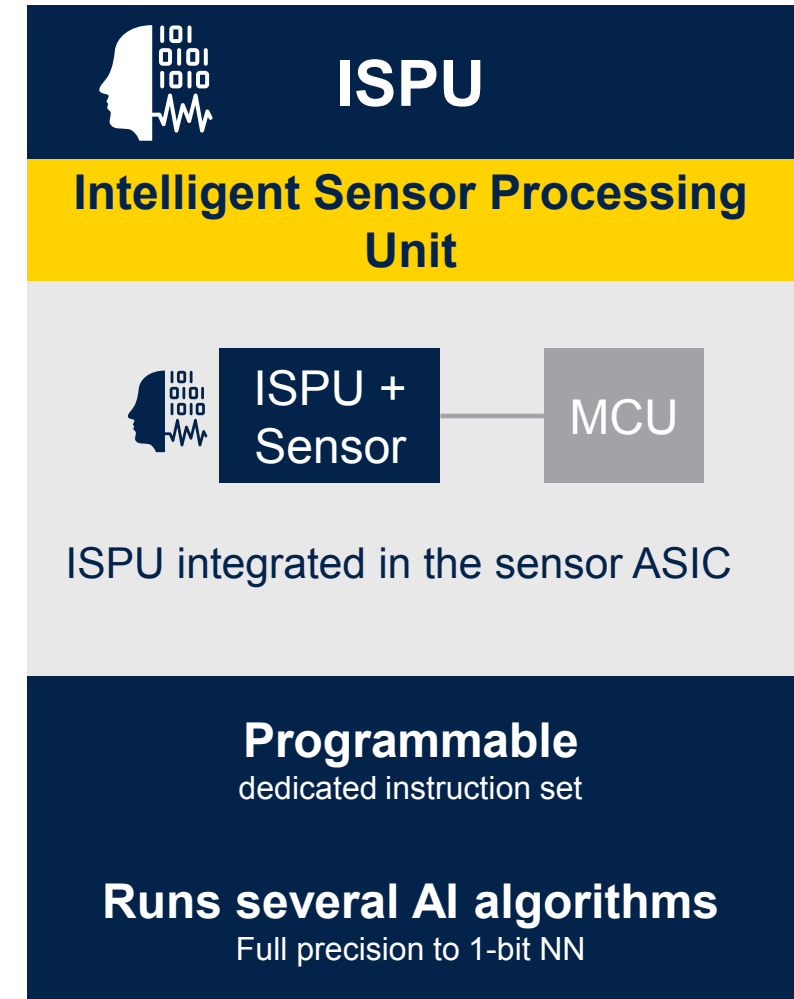| Offline era | Online era | Onlife era |
|---|---|---|
| Fragmented | Connected | Trained |

Local efficiency

Global efficiency

SUSTAINABLE TECHNOLOGY
by ST

life.augmented

# From "on the Edge" to "in the Edge"

| Sensor + MCU | rPU | ISPU |
|---|---|---|
| **Microcontroller** | **reconfigurable Processing Unit** | **Intelligent Sensor Processing Unit** |
| Sensor — MCU | rPU+ Sensor — MCU | ISPU + Sensor — MCU |
| MCU standalone or hosted in the sensor package | rPU integrated in the sensor ASIC | ISPU integrated in the sensor ASIC |
| **Standard** MCU runs the algorithms | **Optimized** reconfigured through register setting | **Programmable** dedicated instruction set |
| **Runs any kind of SW** provided it fits the MCU specs | **Constrained** runs same model/mapping | **Runs several AI algorithms** Full precision to 1-bit NN |

# One solution cannot fit all, but ISPU comes close

| Sensor + MCU | | rPU | | ISPU | |
|---|---|---|---|---|---|
| Sensor — MCU | | rPU+ Sensor — MCU | | ISPU + Sensor — MCU | |

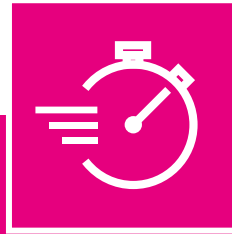| Parameters | | Benefits | | Benefits | | Benefits |
|---|---|---|---|---|---|---|
| **Flexibility** | ▪▪▪▪▪ | ad-hoc coding | ▪□□□□ | | ▪▪▪▪□ | Optimized for inertial data |
| **Low Power** | ▪▪□□□ | | ▪▪▪▪▪ | rPU adds few uA | ▪▪▪▪□ | Integrated computing cell, MCU in standby with sensor wakeup |
| **Data transfer optimization Sensor - MCU** | ▪□□□□ | | ▪▪▪□□ | rPU sends pre-processed data | ▪▪▪▪▪ | Intelligent local processing |

21

# The 3 design criteria for working in the Edge

## Local

In the Edge: data privacy, low power…

## Fast

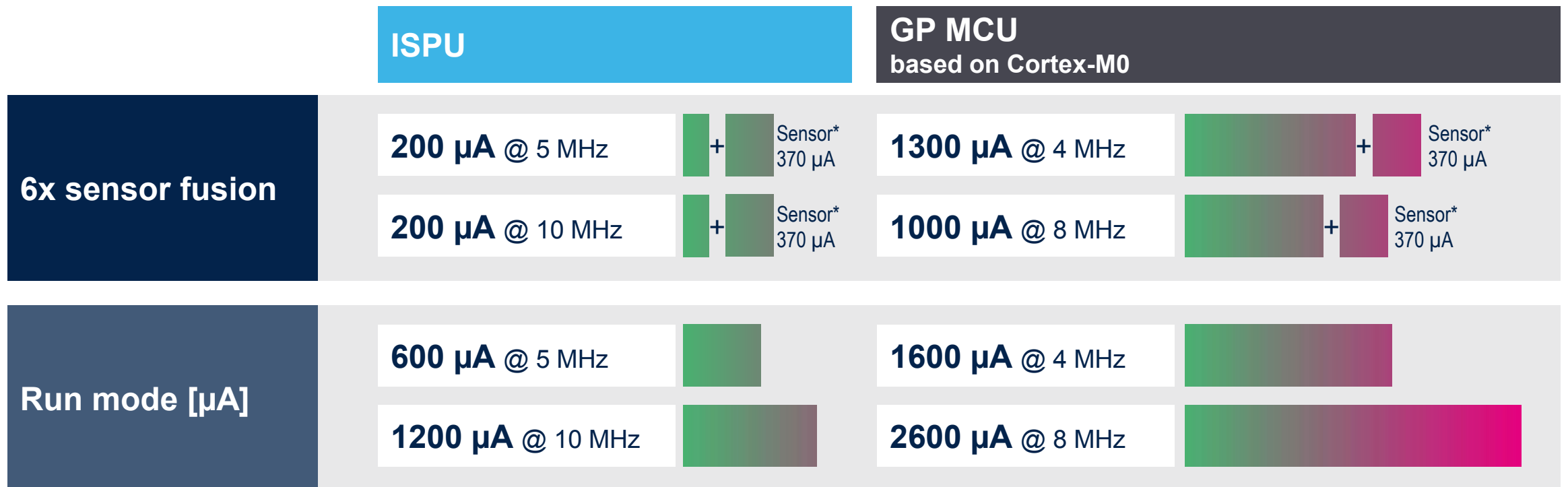Ad-hoc processor customization for real-time execution

## Intelligent

Runs complex AI analyses and takes actions

*life.augmented*
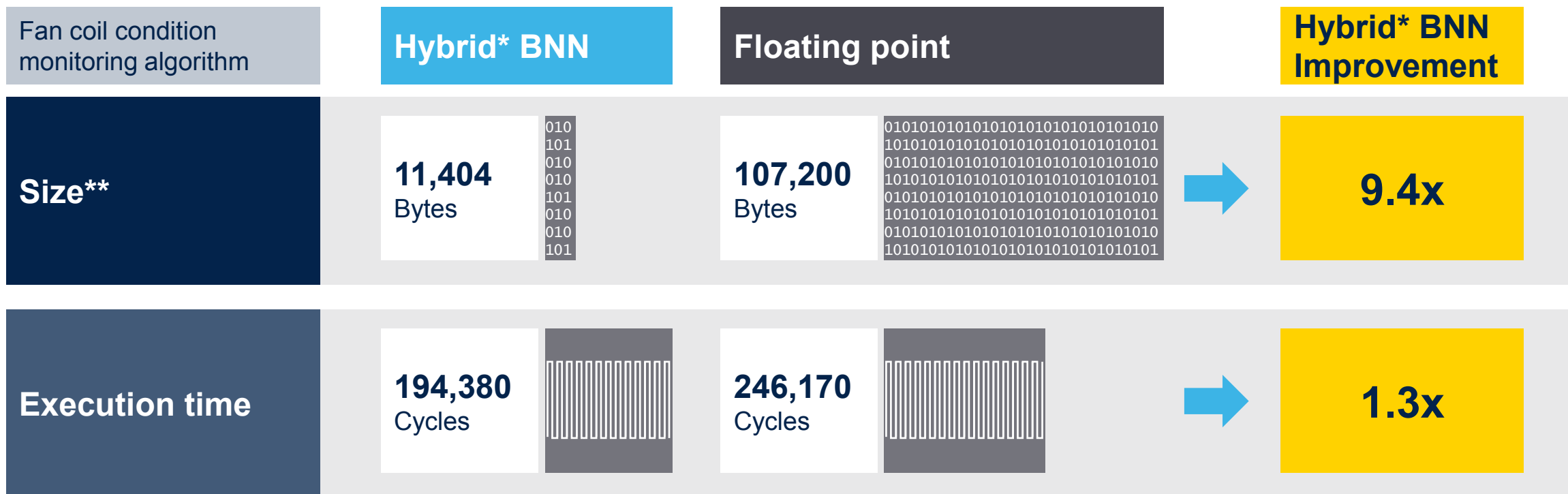
# Sensor fusion on ISPU consumes far less current

**5x less current consumption for sensor fusion on ISPU than on M0**
**Below 600 µA for sensor fusion in the edge**

| | ISPU | | GP MCU based on Cortex-M0 | |
|---|---|---|---|---|
| **6x sensor fusion** | **200 µA** @ 5 MHz | + Sensor* 370 µA | **1300 µA** @ 4 MHz | + Sensor* 370 µA |
| | **200 µA** @ 10 MHz | + Sensor* 370 µA | **1000 µA** @ 8 MHz | + Sensor* 370 µA |
| **Run mode [µA]** | **600 µA** @ 5 MHz | | **1600 µA** @ 4 MHz | |
| | **1200 µA** @ 10 MHz | | **2600 µA** @ 8 MHz | |

*Accelerometer + Gyroscope low-power mode @ ODR 104 Hz

# Running Hybrid Binary Neural Network (BNN) for condition monitoring

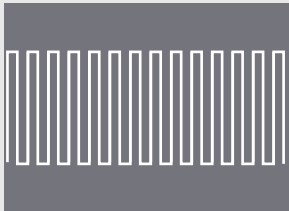**ISPU makes solutions ready for Onlife with faster and smaller algos**

| Fan coil condition monitoring algorithm | Hybrid* BNN | Floating point | Hybrid* BNN Improvement |
|---|---|---|---|
| **Size**** | 11,404 Bytes | 107,200 Bytes | **9.4x** |
| **Execution time** | 194,380 Cycles | 246,170 Cycles | **1.3x** |

\* Some layers are floating point activations with binary weights, some are fully binarized (weights and activations). ISPU set at 5MHz frequency

\*\* Full Application size: data + code + internal buffers + system libs

# BNN on ISPU far outperforms floating point

**BNN on ISPU delivers over 10x better performance than floating point ISPU can now run dense SW layers in the Edge**

| Microbenchmark* [single dense layer 128x64] | Full BNN | Floating point | BNN Improvement |
|---|---|---|---|
| **Size**** | **10,372** Bytes | **106,324** Bytes | **10.3x** |
| **Execution time** | **4,366** Cycles | **74,210** Cycles | **17.0x** |

\*  Kernel size = 128 and number of kernels = 64. ISPU set at 5MHz frequency
\*\*  Full Application size: data + code + internal buffers + system libs

# In-sensor Machine Learning & Deep Learning

**ST ISPU delivers more options and greater freedom**

## Compilation Tool

- Compiler (GNU) / Assembler (GNU) / Linker (GNU)
- Neural network library generation from high level tools (Keras, Tensorflow, etc.)
- Ad hoc optimization for ISPU target

## IDE Tools

- Source-level debugger (GNU)  / On-chip debugger
- Simulator (STMicroelectronics)
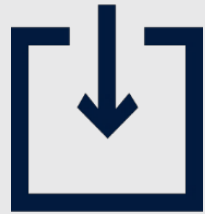- Eclipse graphical interface

## Runtime

- Platform SDK / Peripheral drivers
- Platform libraries

# ISPU with NanoEdge™ AI for self-learning solutions

by

**NANOEDGE AI**

**Onlife-ready:**
classify data patterns and detect in the edge

**Commercial libraries**
ready to deploy on ISPU

**Reference design**
with customization and support

**Industrial IoT**

**Personal Electronics**

Distance
5.73 km

Duration
0:27:34 h

Average HR
130 bpm

Pace
12.45 km/h

Calories burned
280 cal

# Sensor's semiconductors challenges

**AI @ Edge**

ADVANCED TOOLS FOR INCREASED PRODUCTIVITY

ECOSYSTEM COMPATIBILITY WITH AI TOOLS

PROVEN TECHNOLOGY ARCHITECTURE IN SUPER TINY PACKAGE

LIMITED LOGIC AND MEMORY STORAGE FOR EDGE AI

OPTIMAL POWER (uW) PER PERFORMANCE CAPABILITIES

# tinyML Summit 2022 Sponsors

# Copyright Notice

## www.tinyML.org