

# tinyML<sup>®</sup> Talks

*Enabling Ultra-low Power Machine Learning at the Edge*

“Better productivity leveraging AI community driven interoperability”

Danilo Pietro Pau- STMicroelectronics

July 20, 2021



[www.tinyML.org](http://www.tinyML.org)

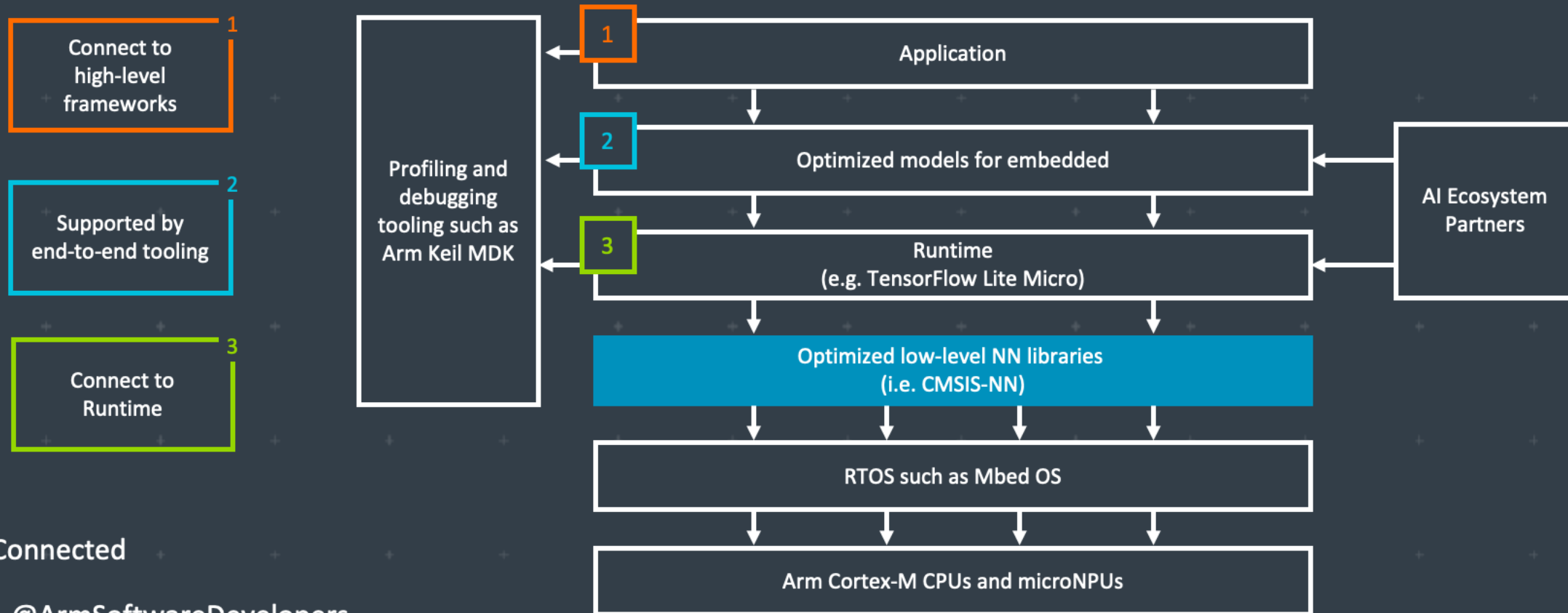


# tinyML Talks Sponsors



Additional Sponsorships available – contact [Olga@tinyML.org](mailto:Olga@tinyML.org) for info

# Arm: The Software and Hardware Foundation for tinyML



Stay Connected

 @ArmSoftwareDevelopers

 @ArmSoftwareDev

Resources: [developer.arm.com/solutions/machine-learning-on-arm](https://developer.arm.com/solutions/machine-learning-on-arm)



# WE USE AI TO MAKE OTHER AI FASTER, SMALLER AND MORE POWER EFFICIENT



**Automatically compress** SOTA models like MobileNet to <200KB with **little to no drop in accuracy** for inference on resource-limited MCUs



**Reduce** model optimization trial & error from weeks to days using Deeplite's **design space exploration**



**Deploy more** models to your device without sacrificing performance or battery life with our **easy-to-use software**

BECOME BETA USER [bit.ly/testdeeplite](https://bit.ly/testdeeplite)



# TinyML for all developers



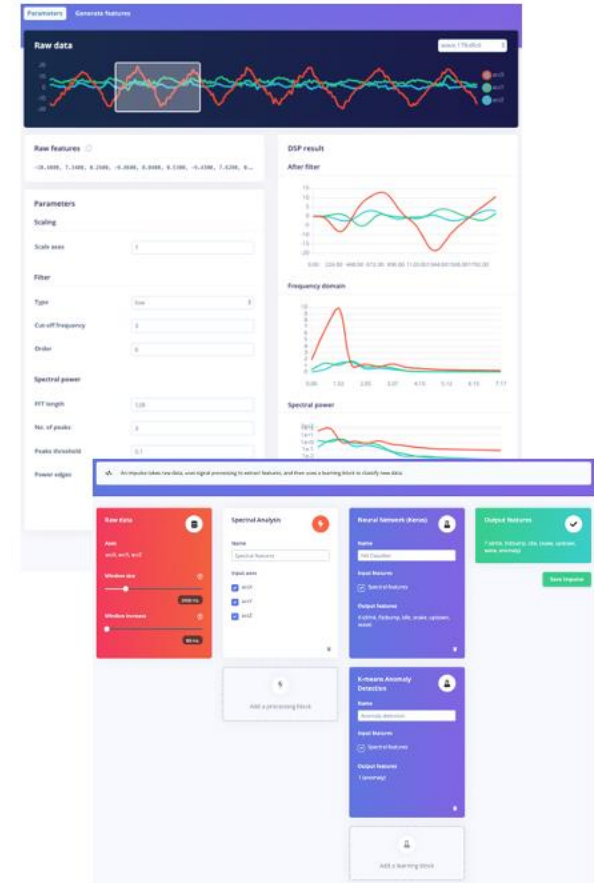
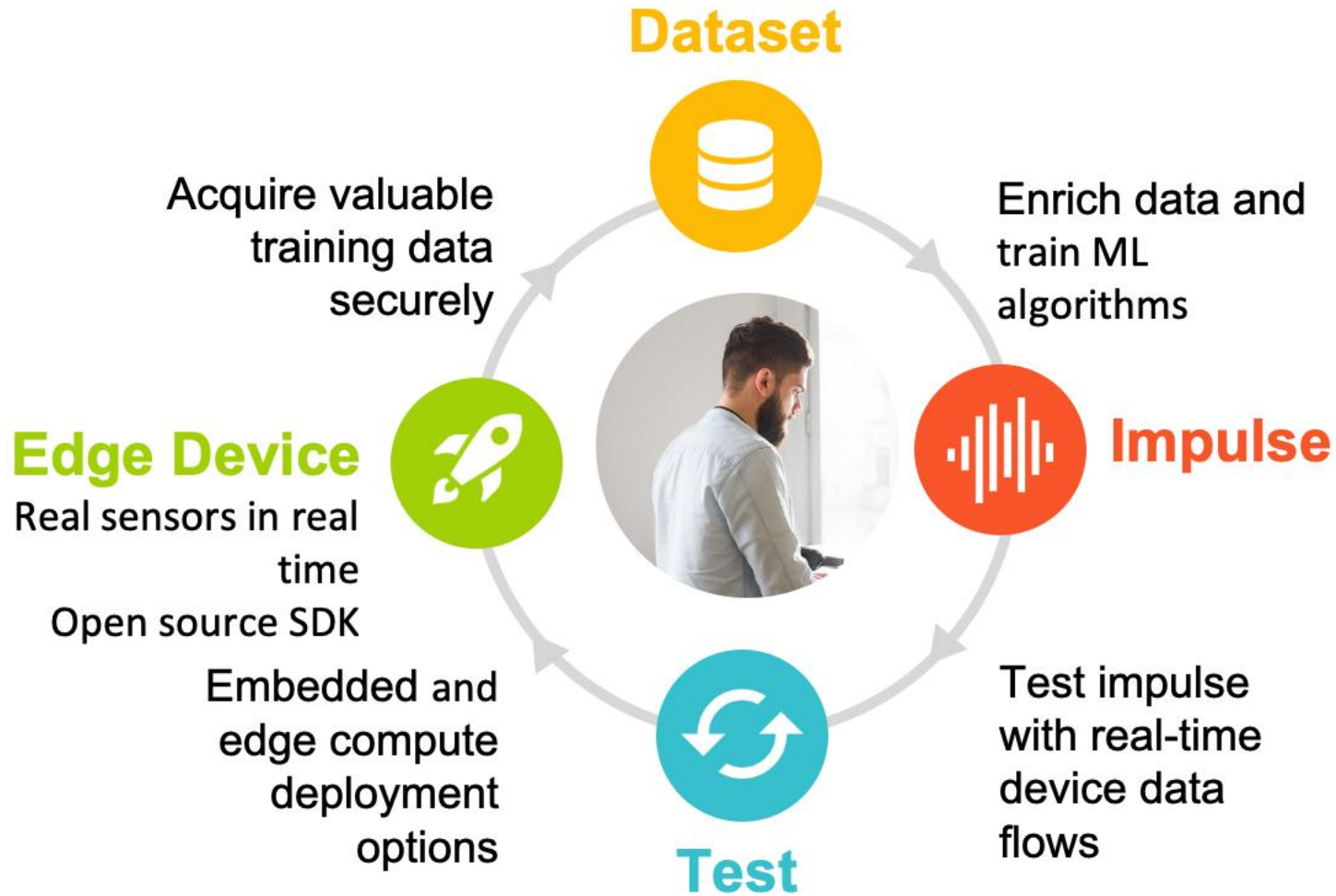
C++ library



Arduino library



WebAssembly

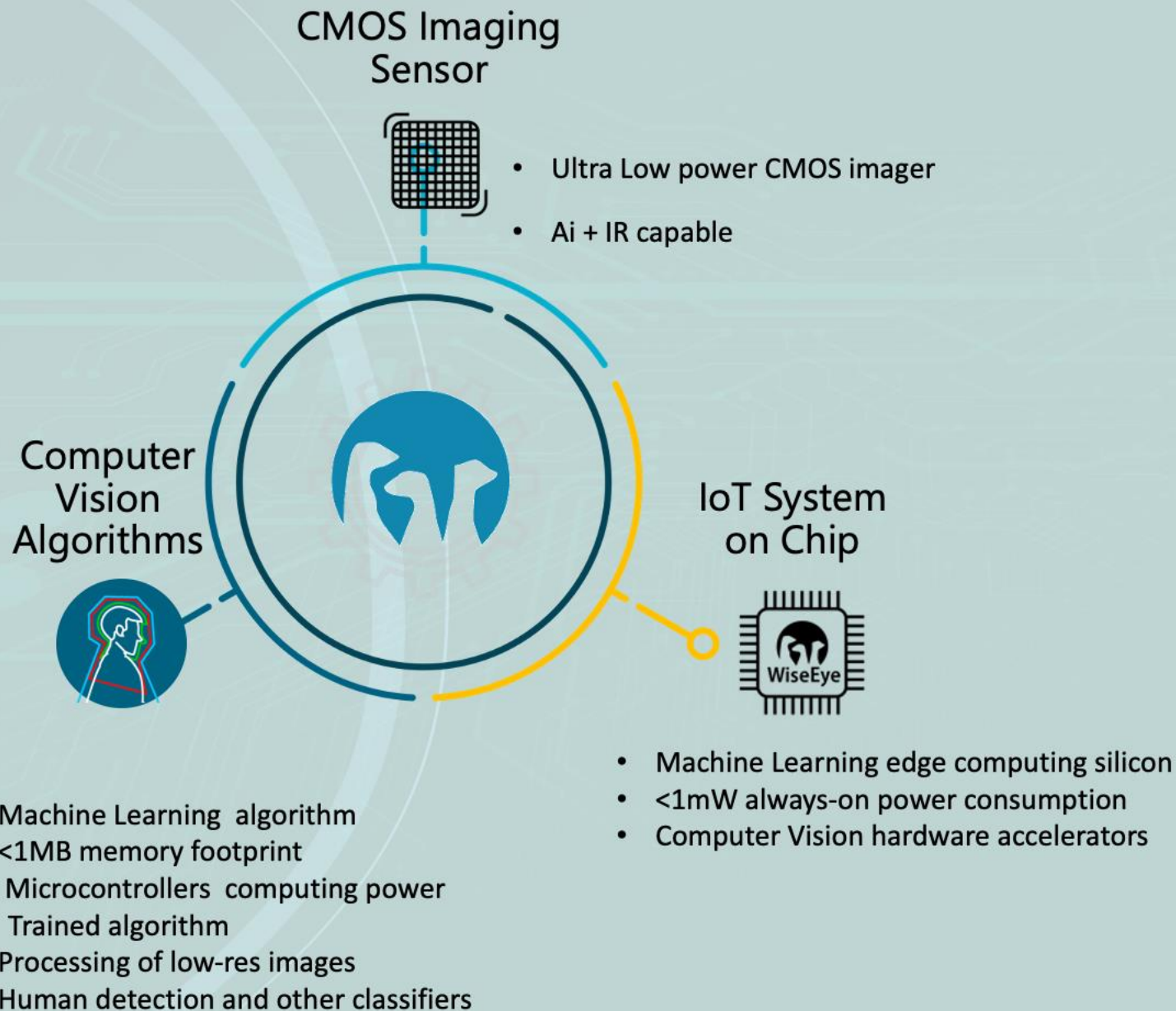


[www.edgeimpulse.com](http://www.edgeimpulse.com)



# The Eye in IoT

## Edge AI Visual Sensors



[info@emza-vs.com](mailto:info@emza-vs.com)



# Enabling the next generation of **Sensor and Hearable products** to process rich data with energy efficiency

Visible Image



Sound



IR Image



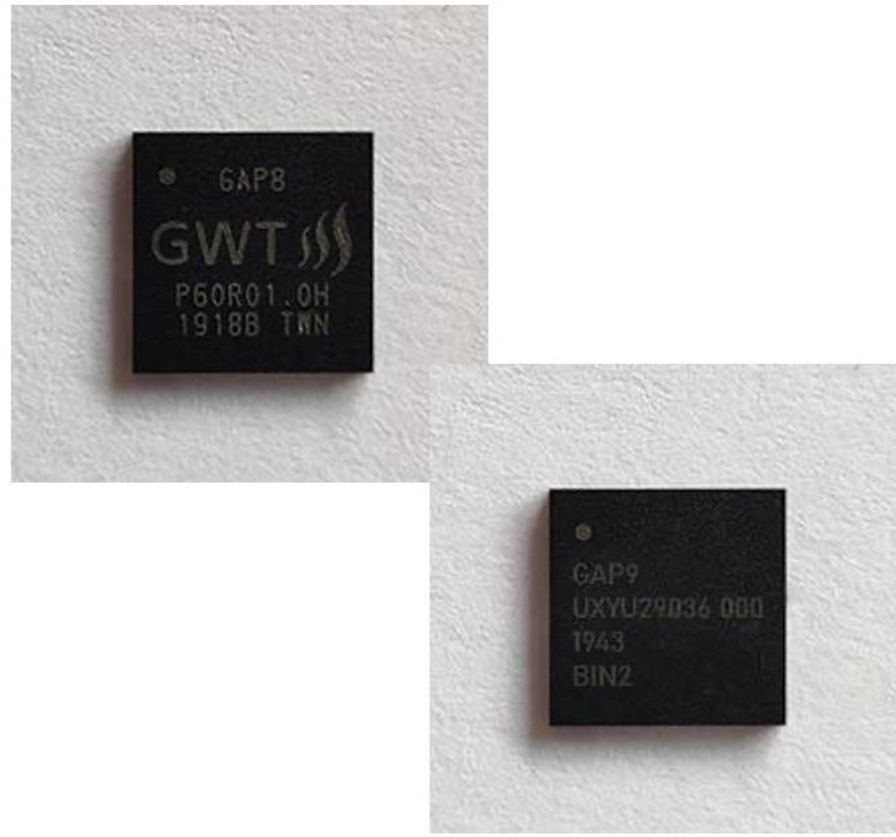
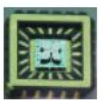
Radar



Bio-sensor



Gyro/Accel



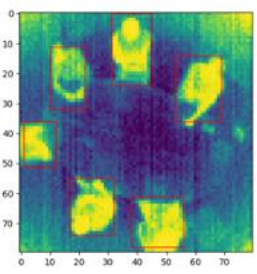
Wearables / Hearables



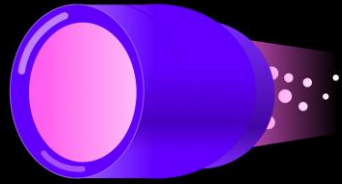
Battery-powered consumer electronics



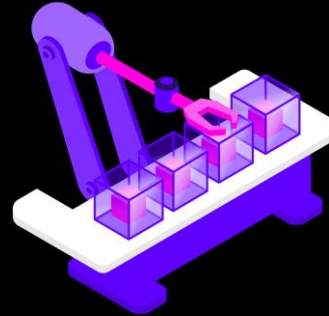
IoT Sensors



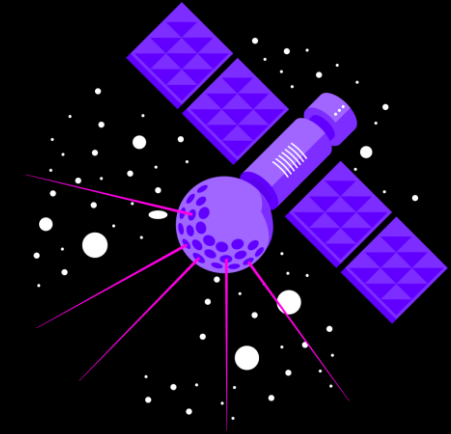
# Distributed infrastructure for TinyML apps



**Develop at warp speed**



**Automate deployments**



**Device orchestration**

**HOTG is building the distributed infrastructure to pave the way for AI enabled edge applications**





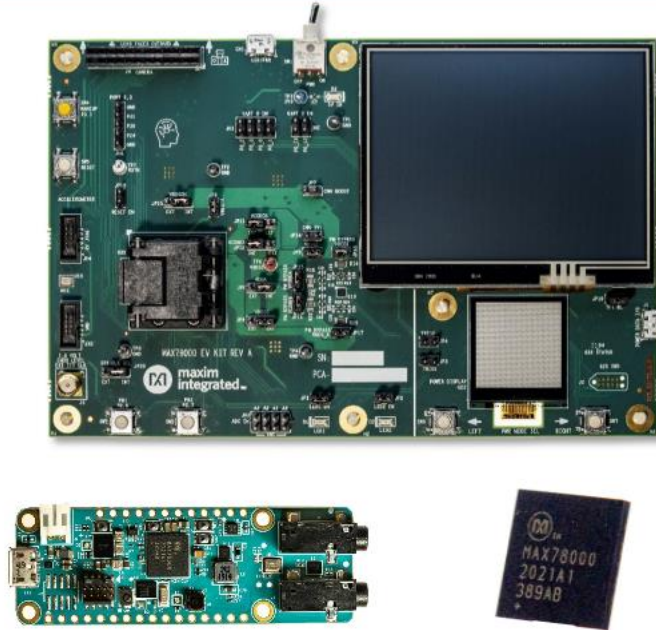
# Latent AI

Adaptive AI for the Intelligent Edge

[latent.ai](https://latent.ai)

## Maxim Integrated: Enabling Edge Intelligence

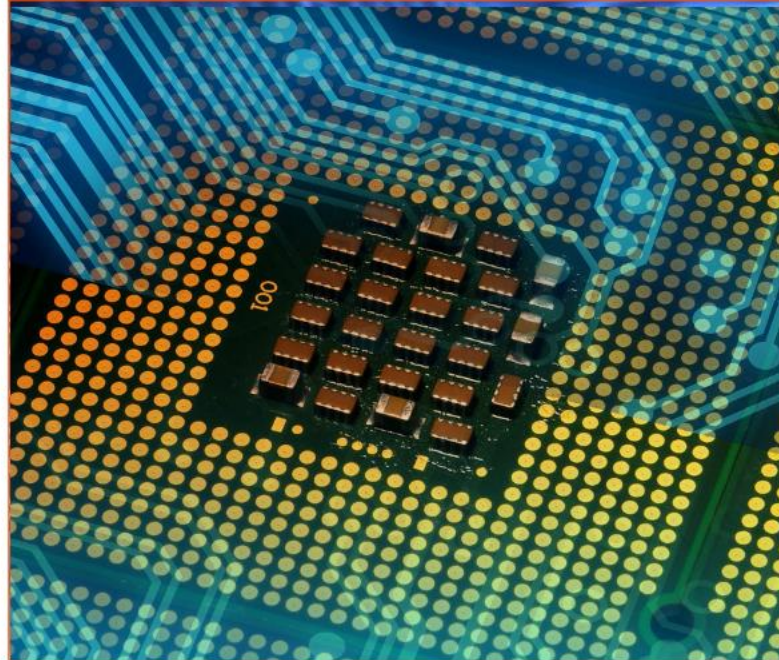
### Advanced AI Acceleration IC



The new MAX78000 implements AI inferences at low energy levels, enabling complex audio and video inferencing to run on small batteries. Now the edge can see and hear like never before.

[www.maximintegrated.com/MAX78000](http://www.maximintegrated.com/MAX78000)

### Low Power Cortex M4 Micros



Large (3MB flash + 1MB SRAM) and small (256KB flash + 96KB SRAM, 1.6mm x 1.6mm) Cortex M4 microcontrollers enable algorithms and neural networks to run at wearable power levels.

[www.maximintegrated.com/microcontrollers](http://www.maximintegrated.com/microcontrollers)

### Sensors and Signal Conditioning



Health sensors measure PPG and ECG signals critical to understanding vital signs. Signal chain products enable measuring even the most sensitive signals.

[www.maximintegrated.com/sensors](http://www.maximintegrated.com/sensors)

# Qeexo AutoML

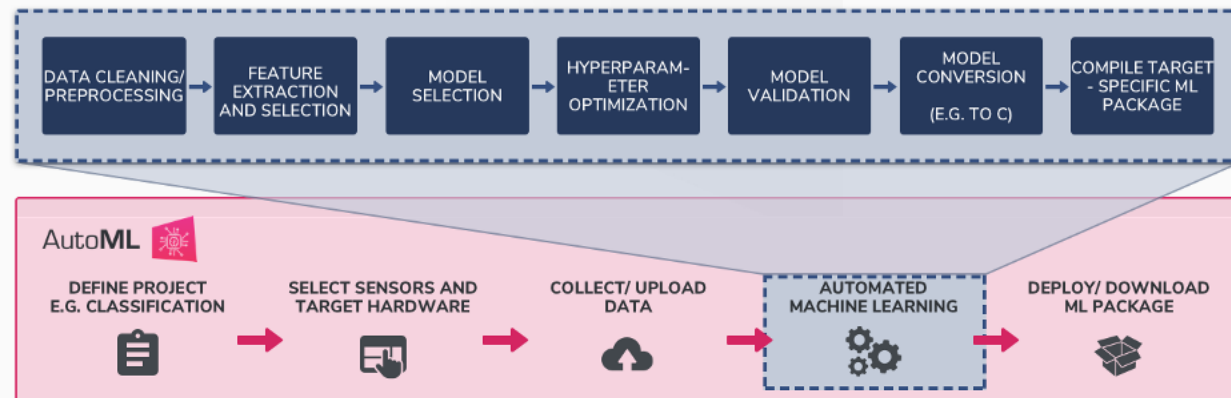


Automated Machine Learning Platform that builds tinyML solutions for the Edge using sensor data

## Key Features

- Supports 17 ML methods:
  - Multi-class algorithms: GBM, XGBoost, Random Forest, Logistic Regression, Gaussian Naive Bayes, Decision Tree, Polynomial SVM, RBF SVM, SVM, CNN, RNN, CRNN, ANN
  - Single-class algorithms: Local Outlier Factor, One Class SVM, One Class Random Forest, Isolation Forest
- Labels, records, validates, and visualizes time-series sensor data
- On-device inference optimized for low latency, low power consumption, and small memory footprint applications
- Supports Arm® Cortex™ - M0 to M4 class MCUs

## End-to-End Machine Learning Platform



For more information, visit: [www.qeexo.com](http://www.qeexo.com)

## Target Markets/Applications

- Industrial Predictive Maintenance
- Smart Home
- Wearables
- Automotive
- Mobile
- IoT



**Qualcomm**  
AI research

# Advancing AI research to make efficient AI ubiquitous

## Power efficiency

Model design, compression, quantization, algorithms, efficient hardware, software tool

## Personalization

Continuous learning, contextual, always-on, privacy-preserved, distributed learning

## Efficient learning

Robust learning through minimal data, unsupervised learning, on-device learning

## A platform to scale AI across the industry



### Perception

Object detection, speech recognition, contextual fusion



### Reasoning

Scene understanding, language understanding, behavior prediction



### Action

Reinforcement learning for decision making



Edge cloud



Cloud



IoT/IloT



Automotive



Mobile





# Reality AI<sup>®</sup>

## Add Advanced Sensing to your Product with Edge AI / TinyML

<https://reality.ai>



[info@reality.ai](mailto:info@reality.ai)



[@SensorAI](https://twitter.com/SensorAI)



[Reality AI](https://www.linkedin.com/company/reality-ai)

## Pre-built Edge AI sensing modules, plus tools to build your own

### Reality AI solutions

Prebuilt sound recognition models for  
indoor and outdoor use cases

Solution for industrial anomaly detection

Pre-built automotive solution that lets cars  
“see with sound”

### Reality AI Tools<sup>®</sup> software

Build prototypes, then turn them into  
real products

Explain ML models and relate the function  
to the physics

Optimize the hardware, including  
sensor selection and placement



# Build Smart IoT Sensor Devices From Data

SensiML pioneered TinyML software tools that auto generate AI code for the intelligent edge.

- End-to-end AI workflow
- Multi-user auto-labeling of time-series data
- Code transparency and customization at each step in the pipeline

We enable the creation of production-grade smart sensor devices.



[sensiml.com](https://sensiml.com)





# SynSense

**SynSense** builds **sensing and inference** hardware for **ultra-low-power** (sub-mW) **embedded, mobile and edge** devices. We design systems for **real-time always-on smart sensing**, for audio, vision, IMUs, bio-signals and more.

<https://SynSense.ai>



# SYNTIANT

[Syntiant Corp.](#) is moving artificial intelligence and machine learning from the cloud to edge devices. Syntiant's chip solutions merge deep learning with semiconductor design to produce ultra-low-power, high performance, deep neural network processors. These network processors enable always-on applications in battery-powered devices, such as smartphones, smart speakers, earbuds, hearing aids, and laptops. Syntiant's Neural Decision Processors™ offer wake word, command word, and event detection in a chip for always-on voice and sensor applications.

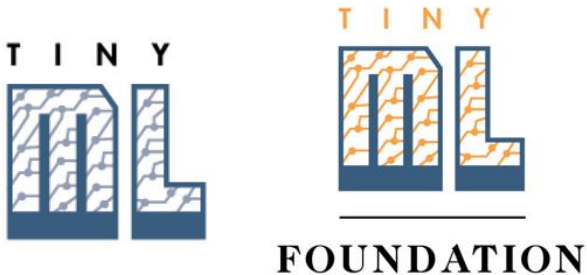
Founded in 2017 and headquartered in Irvine, California, the company is backed by Amazon, Applied Materials, Atlantic Bridge Capital, Bosch, Intel Capital, Microsoft, Motorola, and others. Syntiant was recently named a [CES® 2021 Best of Innovation Awards Honoree](#), [shipped over 10M units worldwide](#), and [unveiled the NDP120](#) part of the NDP10x family of inference engines for low-power applications.

[www.syntiant.com](http://www.syntiant.com)



@Syntiantcorp



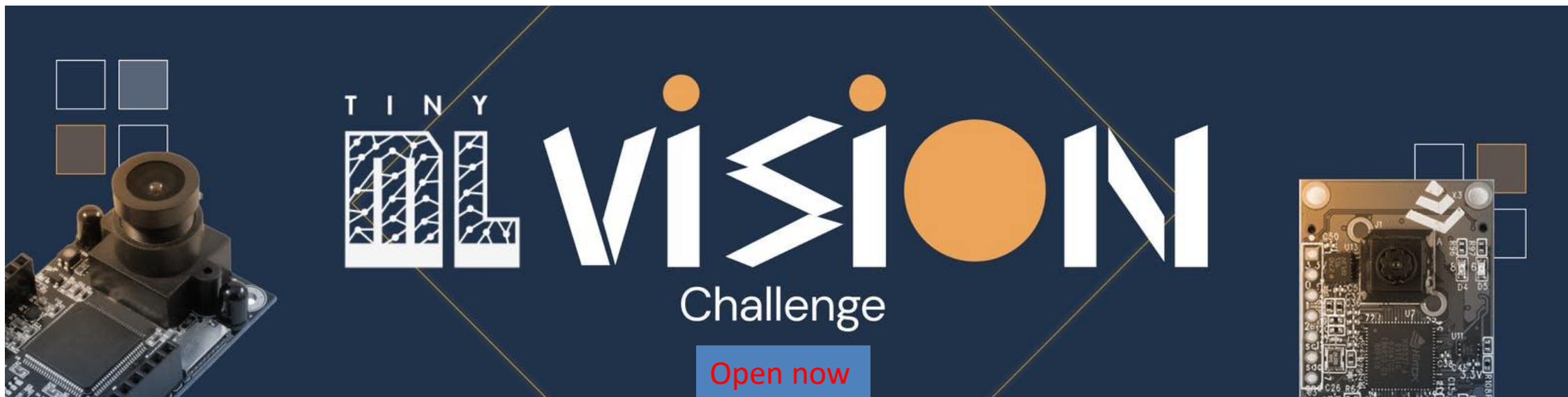


collaboration with



**Focus on:**

(i) developing new use cases/apps for tinyML vision; and (ii) promoting tinyML tech & companies in the developer community



Submissions accepted until August 20<sup>th</sup>, 2021  
Winners announced on September 1, 2021 (\$6k value)  
Sponsorships available: [sponsorships@tinyML.org](mailto:sponsorships@tinyML.org)



<https://www.hackster.io/contests/tinyml-vision>



# Successful tinyML EMEA 2021



- Videos are available on [www.youtube.com/tinyML](https://www.youtube.com/tinyML)

- **4** days of tinyML excitement

- **2** tutorials
- **5** keynotes
- **15** tinyTalks
- **7** lightning talks
- **3** panel discussions & networking
- **16** papers in the Student Forum
- **4** partner sessions
- **16** sponsoring companies

- **58** speakers, **1687** registered attendees!



250 videos with 121k views  
as of July 10, 2021





# Next tinyML Talks

Date	Presenter	Topic / Title
Tuesday, August 3	<b>Vikram Shrivastava,</b> Sr. Director, IoT Marketing, Knowles Corporate	Dedicated Audio Processors at the Edge are the Future of AI

Webcast start time is 8 am Pacific time

Please contact [talks@tinymml.org](mailto:talks@tinymml.org) if you are interested in presenting

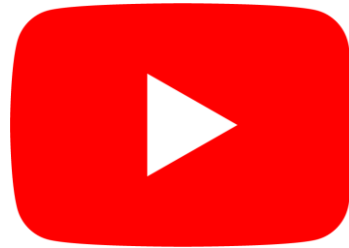


# Reminders

Slides & Videos will be posted tomorrow

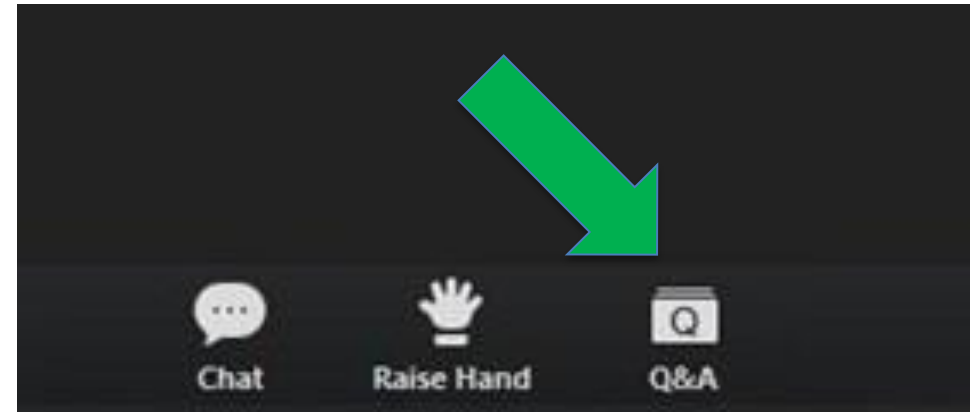


[tinyml.org/forums](https://tinyml.org/forums)



[youtube.com/tinyml](https://youtube.com/tinyml)

Please use the Q&A window for your questions







# Danilo Pau



One year before graduating from the Polytechnic University of Milan in 1992, Danilo PAU joined STMicroelectronics, where he worked on HDMAC and MPEG2 video memory reduction, video coding, embedded graphics, and computer vision. Today, his work focuses on developing solutions for deep learning tools and applications. Since 2019 Danilo has been an IEEE Fellow. Currently serves as a member of IEEE Region 8 Action for Industry and Member of the Machine Learning, Deep Learning and AI in the CE (MDA) Technical Stream Committee IEEE Consumer Electronics Society (CESoc). With over 80 patents, 104 publications, 113 MPEG authored documents and 39 invited talks/seminars at various worldwide Universities and Conferences, Danilo's favorite activity remains mentoring undergraduate students, MSc engineers and PhD students from various universities.



# Better productivity leveraging AI community driven interoperability

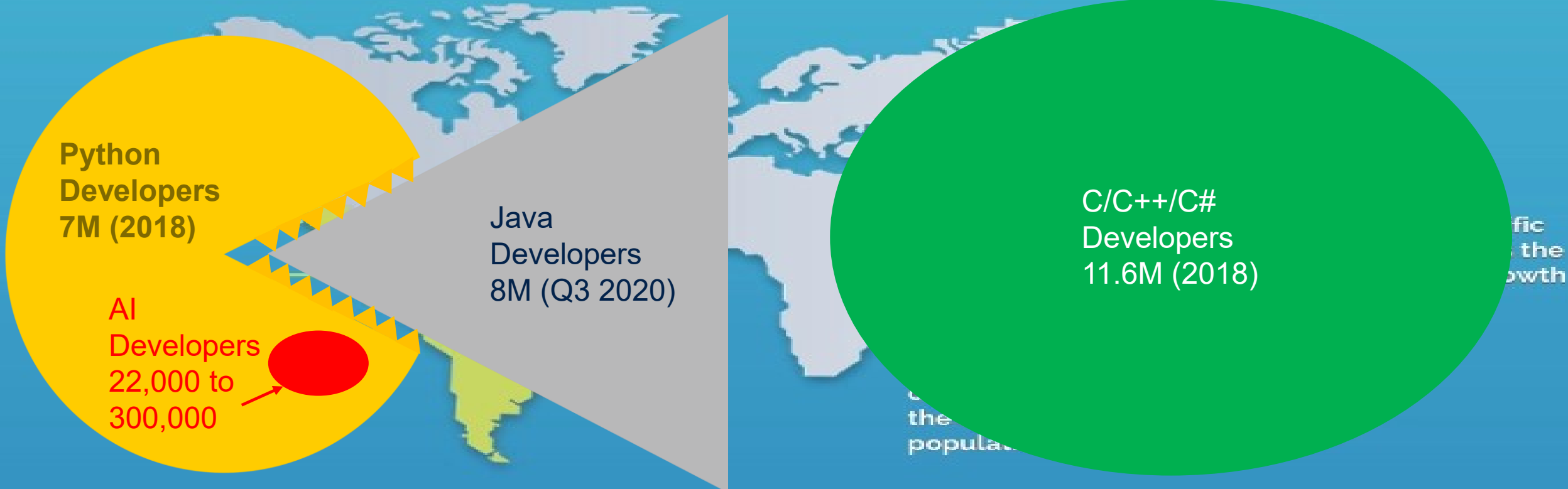
Danilo Pau  
Technical Director, IEEE & ST Fellow  
System Research and Applications  
STMicroelectronics

July 20, 2021

# Global Developer Population and Demographic Study 2019, Vol 1

Source <https://www.daxx.com/blog/development-trends/number-software-developers-world>

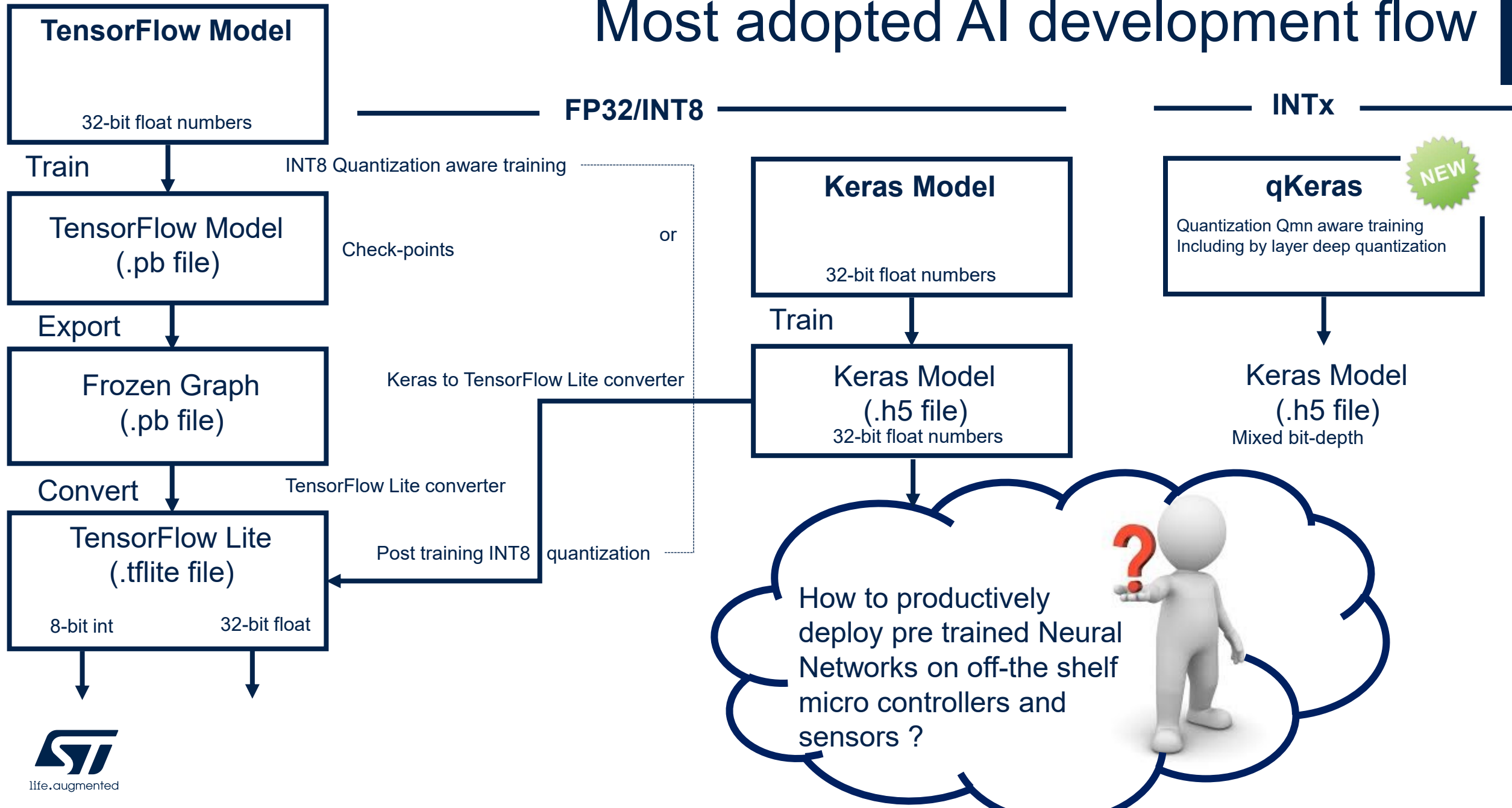
\*<https://www.stateofai2019.com/chapter-6-the-war-for-talent/#:~:text=Estimates%20of%20the%20number%20of,AI%20originated%20in%20academia.>



 **2019: 23.9 million developers**

 **2024: 28.7 million developers**

# Most adopted AI development flow





# Conditions for Deep Learning on MCUs

0  $\text{Accuracy}_{\text{neural networks}} > \text{Accuracy}_{\text{traditional machine learning or hand-crafted algorithm}}$

1  $\text{RAM}_{\text{MCU}} > \text{Activations}_{\text{Network}}$

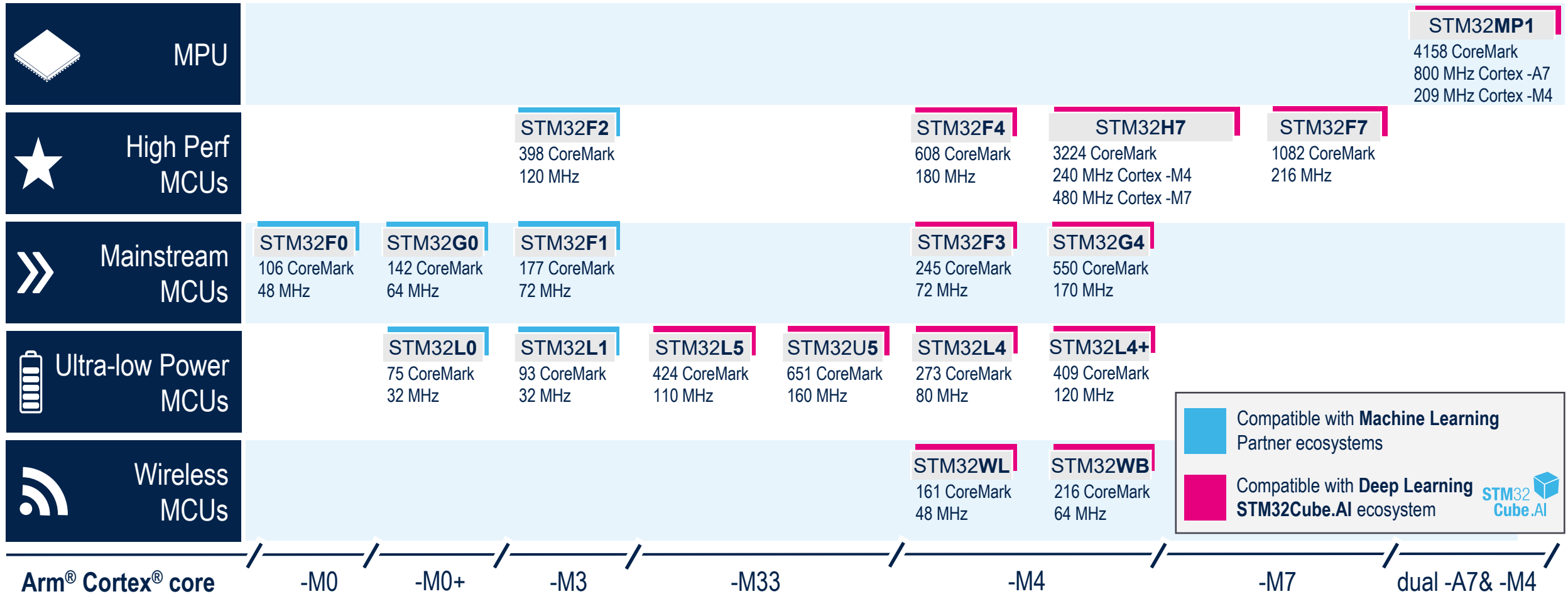
2  $\text{ROM}_{\text{MCU}} > \text{Weights}_{\text{Network}}$

3  $\frac{\text{NN complexity}_{\text{FLOPS}}}{\text{MCU}_{\text{FLOPS/s}}} < \frac{\text{samples in a window}}{f_{\text{sensor sampling rate/s}}}$



# Making AI Accessible on STM32

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



Arm® Cortex® core

-M0

-M0+

-M3

-M33

-M4

-M7

dual -A7 & -M4

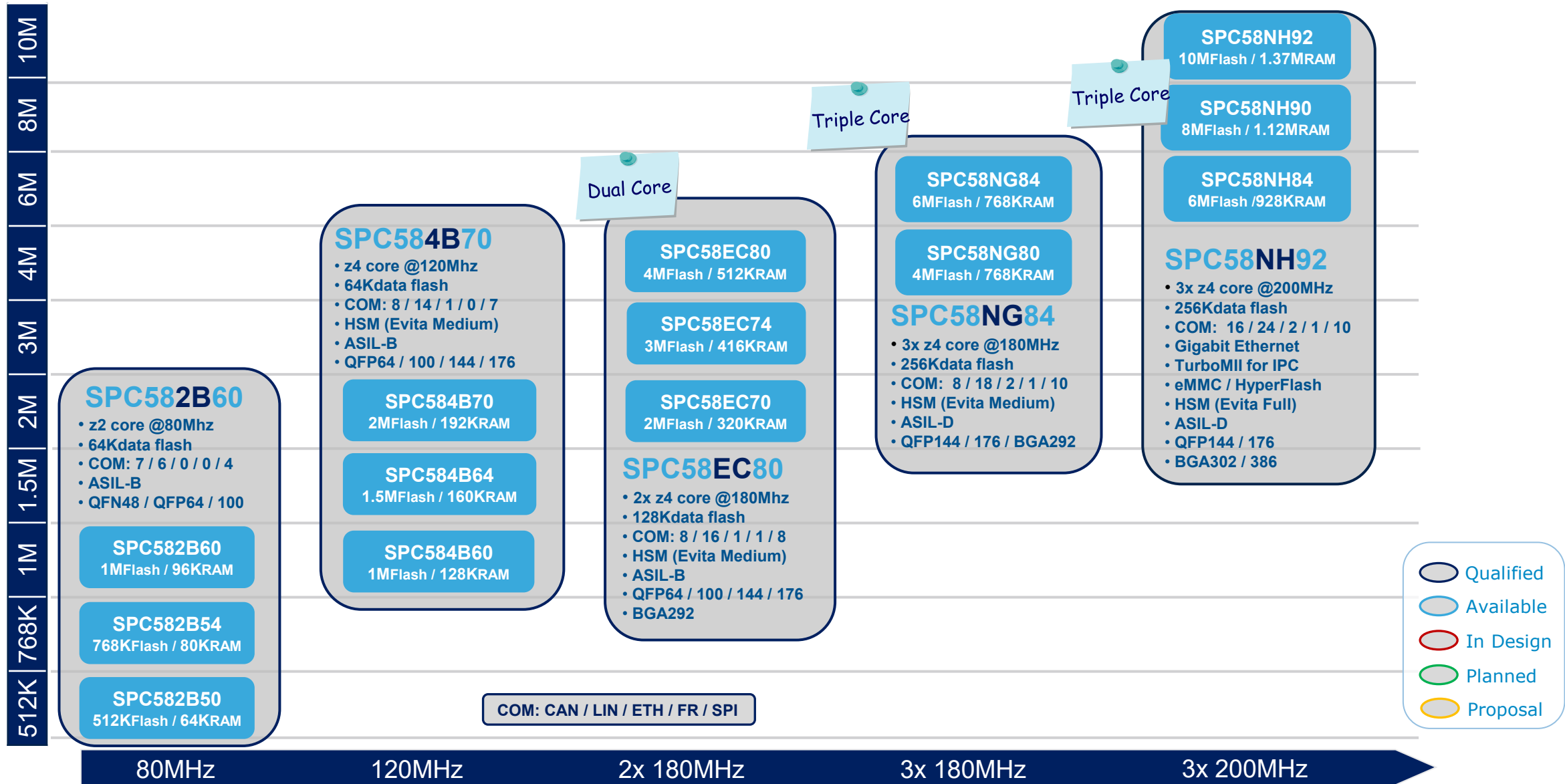
More than 40,000 customers

Over 4 Billion STM32 shipped since 2007





# and on SPC58 - Chorus Family



# How to bridge the AI and embedded communities?



2019: 23.9 million developers



2024: 28.7 million developers





**The message from the embedded developer's community?**

**How to design and deploy resource constrained AI productively ?**

# Key steps for Supervised Deep Learning



1



Clean, label data  
Build NN topology

2



3



Convert NN into  
optimized code for MCU

4



5

Deploy application  
on the field

# Interoperability

Pre-trained Neural Network models  
Deep Learning framework dependent

# Interoperability met across levels !

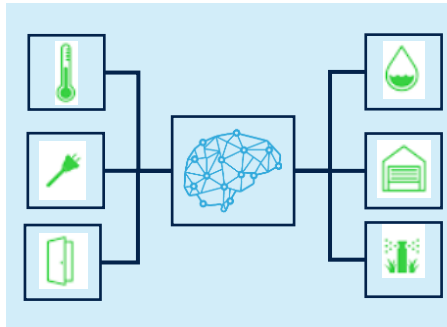
## Everybody else



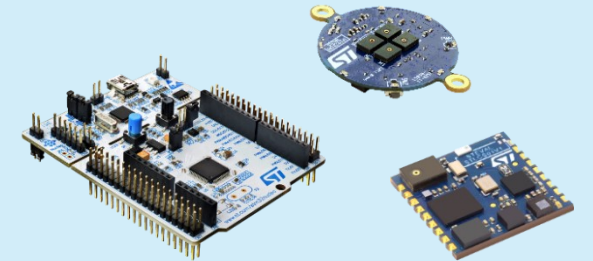
## Google



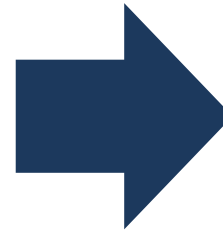
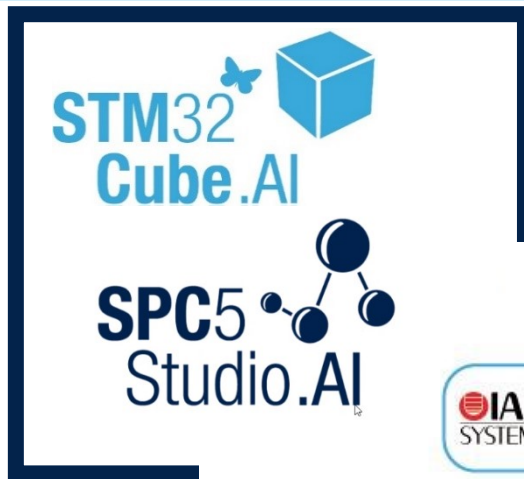
## AI Hub Multiple Neural Networks



Deep Learning SW Solution



## Sensors and OS Agnostic

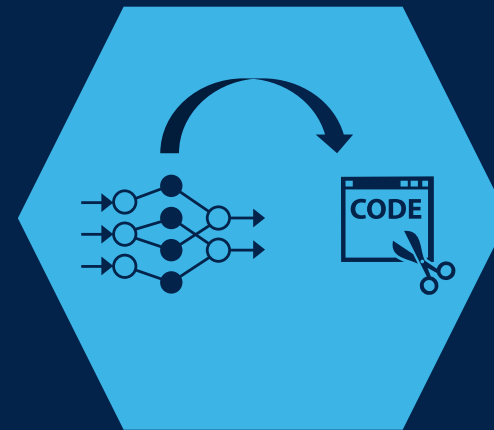


Choose your IDE  
Compiler and Debugger  
Framework Independent

# STM32Cube.AI

STM32   
Cube.AI

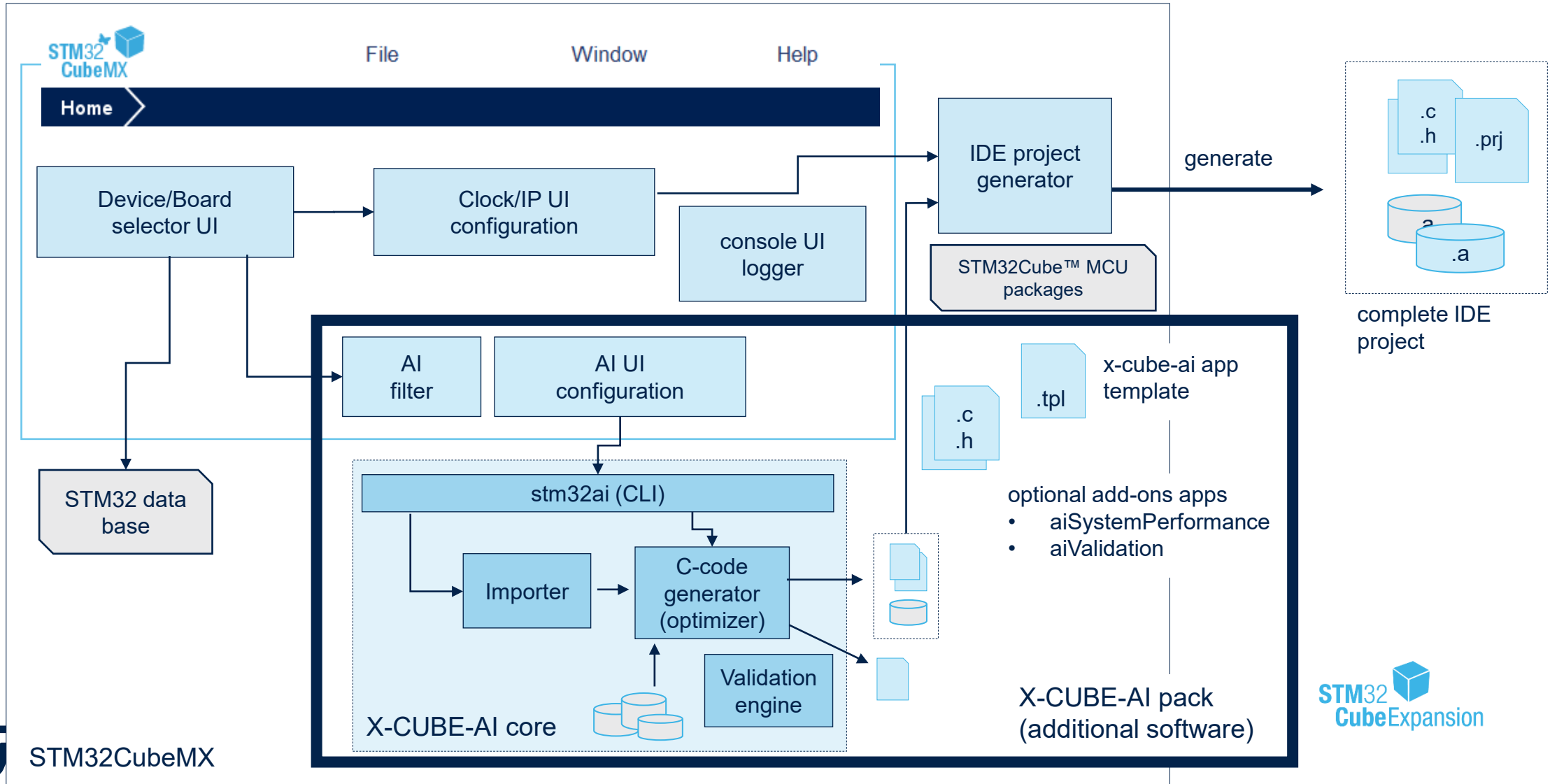
4



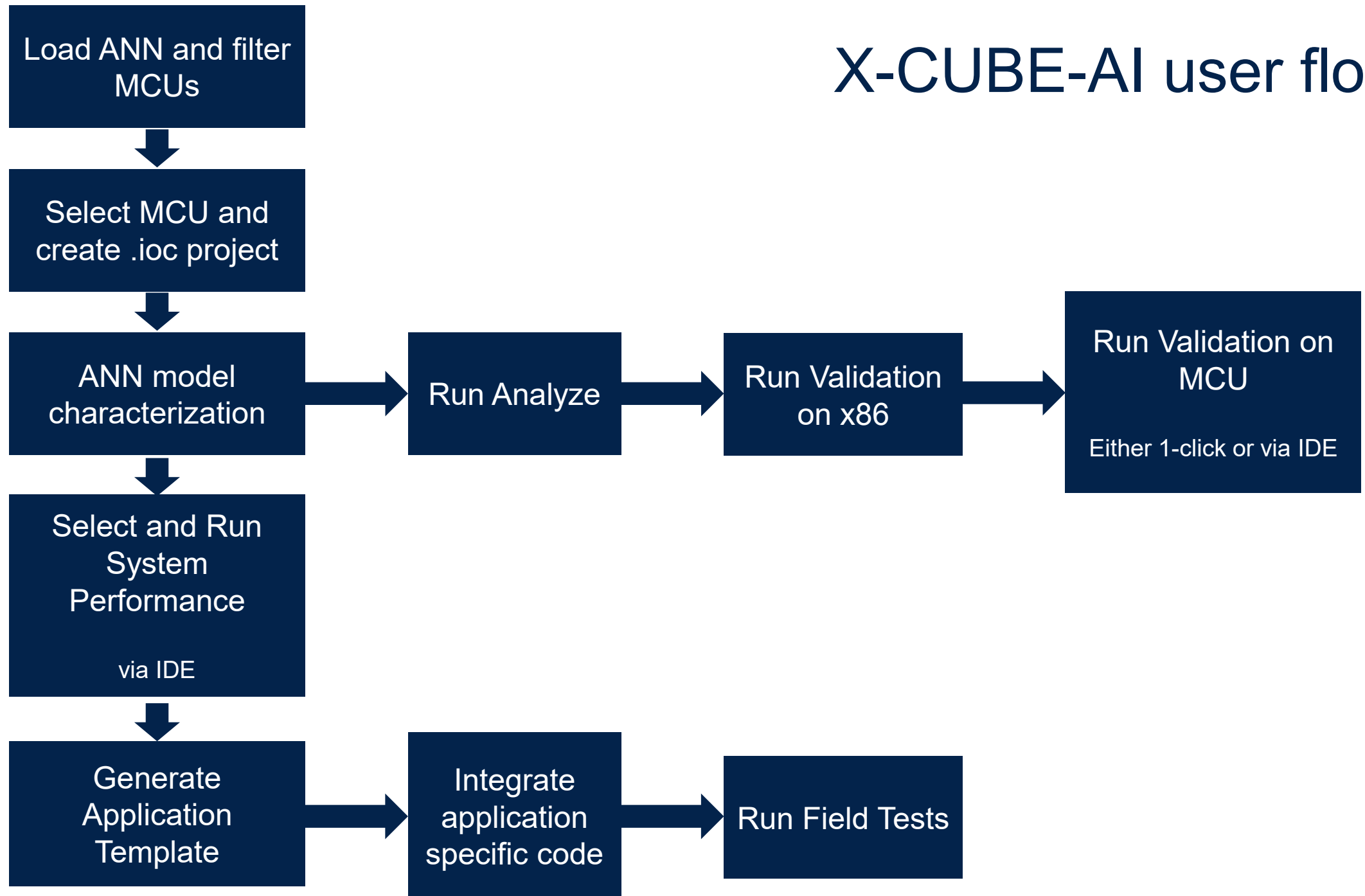
Convert NN into  
optimized code for  
execution



# X-CUBE-AI package a STM32CubeMX additional sw



# X-CUBE-AI user flow



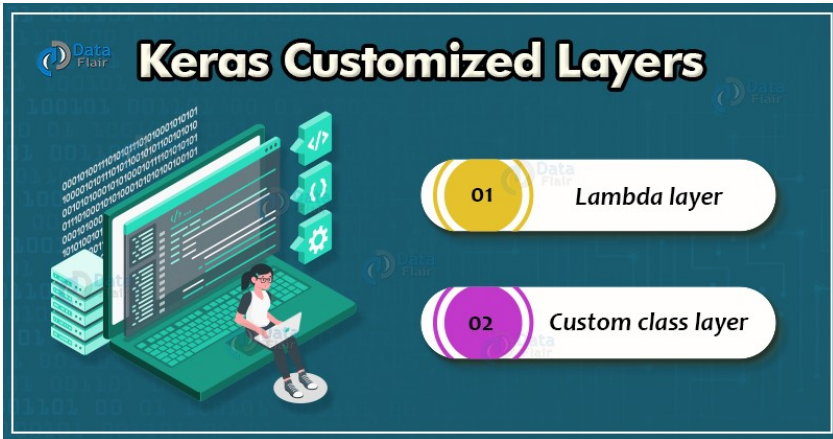


C:\Users\daniilo\_pau\OneDrive - STMicroelectronics\ai>stm32ai\_v700 supported-ops  
Neural Network Tools for STM32AI v1.5.1 (STM.ai v7.0.0-RC8)

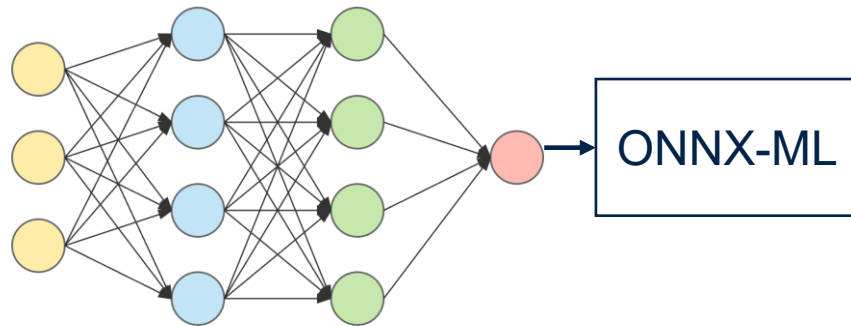
246 operators found

Abs (ONNX), ABS (TFLITE), Acos (ONNX), Acosh (ONNX), Activation (KERAS), ActivityRegularization (KERAS), Add (KERAS), Add (ONNX), ADD (TFLITE), AlphaDropout (KERAS), And (ONNX), ARG\_MAX (TFLITE), ARG\_MIN (TFLITE), ArgMax (ONNX), ArgMin (ONNX), ArrayFeatureExtractor (ONNX), Asin (ONNX), Asinh (ONNX), Atan (ONNX), Atanh (ONNX), Average (KERAS), AVERAGE\_POOL\_2D (TFLITE), AveragePool (ONNX), AveragePooling1D (KERAS), AveragePooling2D (KERAS), BATCH\_TO\_SPACE\_ND (TFLITE), BatchNormalization (KERAS), BatchNormalization (ONNX), Bidirectional (KERAS), Cast (ONNX), CAST (TFLITE), Ceil (ONNX), CEIL (TFLITE), Clip (ONNX), Concat (ONNX), Concatenate (KERAS), CONCATENATION (TFLITE), Constant (ONNX), Conv (ONNX), Conv1D (KERAS), Conv2D (KERAS), Conv2DTranspose (KERAS), CONV\_2D (TFLITE), ConvTranspose (ONNX), Cos (ONNX), COS (TFLITE), Cosh (ONNX), Cropping1D (KERAS), Cropping2D (KERAS), CustomFloorDiv (KERAS), CustomFloorMod (KERAS), CustomPow (KERAS), CustomReshape (KERAS), CustomShape (KERAS), CustomUnpack (KERAS), Dense (KERAS), DEPTHWISE\_CONV\_2D (TFLITE), DepthwiseConv2D (KERAS), DEQUANTIZE (TFLITE), DequantizeLinear (ONNX), Div (ONNX), DIV (TFLITE), Dropout (KERAS), Dropout (ONNX), ELU (KERAS), Elu (ONNX), ELU (TFLITE), Equal (ONNX), EQUAL (TFLITE), Erf (ONNX), Exp (ONNX), EXP (TFLITE), EXPAND\_DIMS (TFLITE), Flatten (KERAS), Flatten (ONNX), Floor (ONNX), FLOOR (TFLITE), FLOOR\_DIV (TFLITE), FLOOR\_MOD (TFLITE), FULLY\_CONNECTED (TFLITE), GaussianDropout (KERAS), GaussianNoise (KERAS), Gemm (ONNX), GlobalAveragePool (ONNX), GlobalAveragePooling1D (KERAS), GlobalAveragePooling2D (KERAS), GlobalMaxPool (ONNX), GlobalMaxPooling1D (KERAS), GlobalMaxPooling2D (KERAS), Greater (ONNX), GREATER (TFLITE), GREATER\_EQUAL (TFLITE), GreaterOrEqual (ONNX), GRU (KERAS), HARD\_SWISH (TFLITE), HardMax (ONNX), HardSigmoid (ONNX), Identity (ONNX), InputLayer (KERAS), InstanceNormalization (ONNX), KaldiNormLayer (KERAS), L2\_NORMALIZATION (TFLITE), LabelEncoder (ONNX), LEAKY\_RELU (TFLITE), LeakyReLU (KERAS), LeakyRelu (ONNX), Less (ONNX), LESS (TFLITE), LESS\_EQUAL (TFLITE), LessOrEqual (ONNX), LOCAL\_RESPONSE\_NORMALIZATION (TFLITE), Log (ONNX), LOG (TFLITE), LOG\_SOFTMAX (TFLITE), LOGICAL\_AND (TFLITE), LOGICAL\_NOT (TFLITE), LOGICAL\_OR (TFLITE), LOGISTIC (TFLITE), LogSoftMax (ONNX), LpNormalization (ONNX), LRN (ONNX), LSTM (KERAS), LSTM (ONNX), LSTM (TFLITE), MatMul (ONNX), Max (ONNX), MAX\_POOL\_2D (TFLITE), Maximum (KERAS), MAXIMUM (TFLITE), MaxPool (ONNX), MaxPooling1D (KERAS), MaxPooling2D (KERAS), Mean (ONNX), MEAN (TFLITE), Min (ONNX), Minimum (KERAS), MINIMUM (TFLITE), MIRROR\_PAD (TFLITE), Mul (ONNX), MUL (TFLITE), Multiply (KERAS), Neg (ONNX), NEG (TFLITE), Not (ONNX), Or (ONNX), PACK (TFLITE), Pad (ONNX), PAD (TFLITE), PADV2 (TFLITE), Permute (KERAS), PlaceholderCustomLayer (KERAS), Pow (ONNX), POW (TFLITE), PReLU (KERAS), PRelu (ONNX), PRELU (TFLITE), QLinearConv (ONNX), QLinearMatMul (ONNX), QUANTIZE (TFLITE), QuantizeLinear (ONNX), Reciprocal (ONNX), REDUCE\_ANY (TFLITE), REDUCE\_MAX (TFLITE), REDUCE\_MIN (TFLITE), REDUCE\_PROD (TFLITE), ReduceL1 (ONNX), ReduceL2 (ONNX), ReduceMax (ONNX), ReduceMean (ONNX), ReduceMin (ONNX), ReduceProd (ONNX), ReduceSum (ONNX), ReduceSumSquare (ONNX), ReLU (KERAS), Relu (ONNX), RELU (TFLITE), RELU6 (TFLITE), RELU\_N1\_TO\_1 (TFLITE), RepeatVector (KERAS), Reshape (KERAS), Reshape (ONNX), RESHAPE (TFLITE), Resize (ONNX), RESIZE\_BILINEAR (TFLITE), RESIZE\_NEAREST\_NEIGHBOR (TFLITE), Round (ONNX), ROUND (TFLITE), Rsqrt (ONNX), RSQRT (TFLITE), Selu (ONNX), SeparableConv1D (KERAS), SeparableConv2D (KERAS), Shape (ONNX), SHAPE (TFLITE), Sigmoid (ONNX), Sign (ONNX), Sin (ONNX), SIN (TFLITE), Sinh (ONNX), Slice (ONNX), SLICE (TFLITE), Softmax (KERAS), Softmax (ONNX), SOFTMAX (TFLITE), Softplus (ONNX), Softsign (ONNX), SPACE\_TO\_BATCH\_ND (TFLITE), SpatialDropout1D (KERAS), SpatialDropout2D (KERAS), SPLIT (TFLITE), Sqrt (ONNX), SQRT (TFLITE), SQUARE (TFLITE), Squeeze (ONNX), SQUEEZE (TFLITE), STRIDED\_SLICE (TFLITE), Sub (ONNX), SUB (TFLITE), Subtract (KERAS), Sum (ONNX), SUM (TFLITE), SVMClassifier (ONNX), SVMRegressor (ONNX), Tan (ONNX), Tanh (ONNX), TANH (TFLITE), TensorFlowOpLayer (KERAS), ThresholdedReLU (KERAS), ThresholdedRelu (ONNX), Tile (ONNX), TILE (TFLITE), TimeDistributed (KERAS), Transpose (ONNX), TRANSPOSE (TFLITE), TRANSPOSE\_CONV (TFLITE), TreeEnsembleRegressor (ONNX), UNIDIRECTIONAL\_SEQUENCE\_LSTM (TFLITE), UNPACK (TFLITE), Unsqueeze (ONNX), Upsample (ONNX), UpSampling1D (KERAS), UpSampling2D (KERAS), Xor (ONNX), ZeroPadding1D (KERAS), ZeroPadding2D (KERAS), ZipMap (ONNX)

# Some advanced features

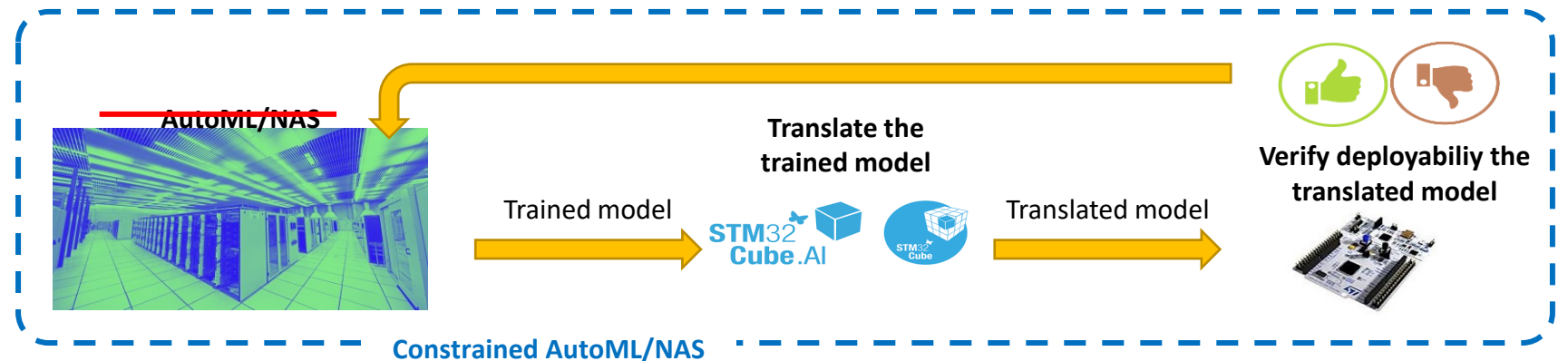


= to express developer creativity, create differentiation, your IPs  
Available since X-CUBE-AI v6



= mix neural and machine learning layers  
Available since X-CUBE-AI v7

Automated deployable design  
Through X-CUBE-AI CLI





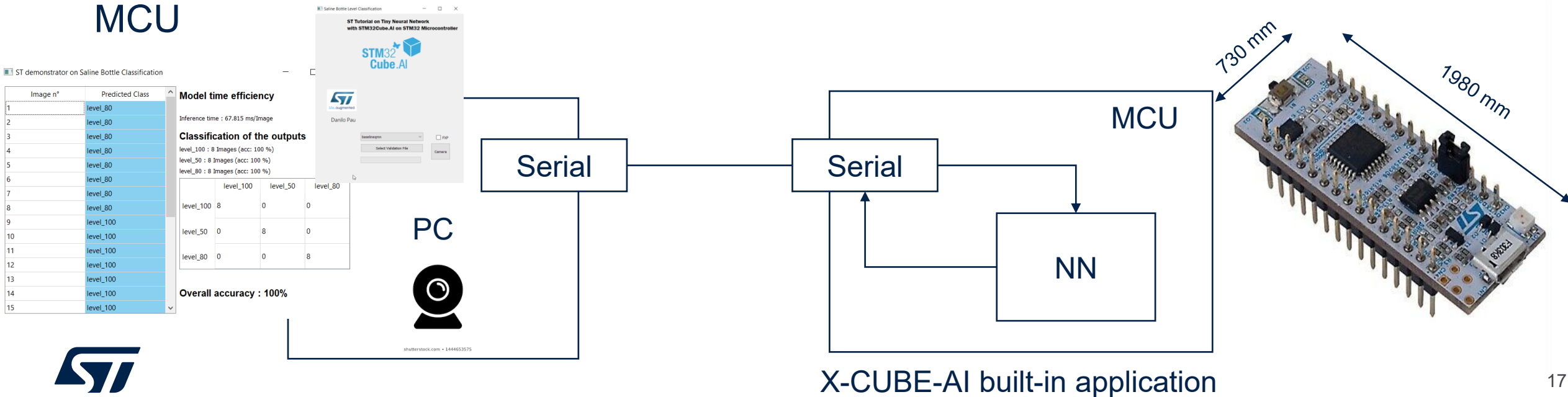
# Saline bottle image classification



	480MHz 2MB FLASH 1MB RAM	MACC	PARAMS	ROM KB	RAM KB	TIME/INFERENC E ms	ACCURACY
1	Baseline FP32	20,863,600	1,247,267	4.76 MB	140	NA	95.3%
2	Baseline FP32 compression=4	20,863,085	1,247,267	1.38 MB	140	212.46	same as 1
3	Baseline INT8	20,776,316	1,247,267	1.19 MB	40.94	69.6	same as 1
4	Reduced FP32	6,707,837	67,891	265.2	48.25	75.29	96%
5	Reduced INT8	6,669,699	67,891	67.11	17.53	35.9	same as 4
6	Tiny1 FP32	451,045	11,547	45.11	17.75	6.51	93.2%
7	Tiny1 INT8	451,051	11,547	11.64	8.45	2.43	same as 6
8	Tiny1DW FP32	253,384	9,419	35.67	18.12	3.84	90.84%
9	Tiny1DW INT8	246,772	9,131	9.46	9.65	1.55	-0.43 vs 8
10	Tiny2 FP32	396,625	7,671	29.96	17.75	6	92.1%
11	Tiny2 INT8	387,687	7,671	7.75	8.37	2.3	+0.43 vs 10
12	Tiny3 Sep FP32	342,749	4,910	19.5	35	6.7	93.5%
13	Tiny3 Sep INT8	342,755	4,910	5.5	9.64	2.63	-0.44 vs 12
14	Tiny4 Sep FP32	43,672	1,126	4.57	12.2	1.4	84.55%
15	Tiny4 Sep INT8	39,023	1,126	1.46	3.16	0.821	same as 14

# Affordable POC for early IoT practitioners

- Attach a MUCLEO-STM32 to the laptop:
  - STM32 MCU runs the model generated with X-CUBE-AI «validation on target» app
  - Image Sensor is attached to PC (webcam), python script reads sensor data and sends to the NUCLEO-STM32, which process and send back results to GUI
  - Communication through a serial port emulated on USB
  - Data information is encoded on PC with the STM32 binary protocol and decoded on MCU



# Affordable POC for early IoT STM32H7 practitioners

The screenshot displays a software interface for a Tiny Neural Network demonstration. The interface is divided into several sections:

- ST demonstrator on Tiny Neural Network: Test Results**: A window showing test results for a single image. It includes a table with the following data:

Image n°	Predicted Class
1	bottle full
- Model time efficiency**: A section indicating the inference time is 3.835 ms/Image.
- Classification of the outputs**: A section showing the classification results for different bottle fill levels:
  - bottle full : 1 Images
  - bottle 80% : 0 Images
  - bottle 50% : 0 Images
- ST Tutorial on Tiny Neural Network with STM32Cube.AI on STM32 Microcontroller**: A main window displaying the ST logo, the STM32Cube.AI logo, and the authors' names: Danilo Pau and Alessandro Carra, V2.0.1. It also features a dropdown menu set to 'tiny1dw' and a 'Refresh NN and camera' button.
- Capturing**: A window showing a live camera feed of a bottle, which is the subject of the neural network's classification.

# Affordable POC for early IoT STM32F3 practitioners

The screenshot displays the STM32Cube.AI software interface. A window titled "ST demonstrator on Tiny Neural Network: Test Results" is open, showing the following data:

Image n°	Predicted Class
1	bottle 50%

Below the table, the text "Inference time : 325.722 ms/Image" is displayed. Under the heading "Classification of the outputs", the following results are listed:

- bottle full : 0 Images
- bottle 80% : 0 Images
- bottle 50% : 1 Images

The background shows the main STM32Cube.AI interface with the title "Tutorial on Tiny Neural Network STM32Cube.AI on STM32 Microcontroller" and the authors "Danilo Pau" and "Alessandro Carra V2.0.1". A "Capturing" window in the foreground shows a video frame of a glass bottle partially filled with liquid.



# (Tiny) Function Packs

Simple, fast, optimized

STM32   
Cube.AI



# STM32 Embedded Software Packages

Part Number	Manufacturer	Description
<u><a href="#">X-LINUX-AI</a></u>	ST	STM32 MPU OpenSTLinux Expansion Pack for AI computer vision application
<u><a href="#">FP-AI-SENSING1</a></u>	ST	STM32Cube function pack for ultra-low power IoT node with artificial intelligence (AI) application based on audio and motion sensing
<u><a href="#">FP-AI-VISION1</a></u>	ST	STM32Cube function pack for high performance STM32 with artificial intelligence (AI) application for Computer Vision
<u><a href="#">FP-AI-NANOEDG1</a></u>	ST	Artificial Intelligence (AI) condition monitoring function pack for STM32Cube
<u><a href="#">FP-AI-FACEREC</a></u>	ST	Artificial Intelligence (AI) face recognition function pack for STM32Cube
<u><a href="#">FP-AI-CTXAWARE1</a></u>	ST	STM32Cube function pack for ultra-low power context awareness with distributed artificial intelligence (AI)



# Context awareness:

[https://www.youtube.com/watch?v=I\\_XqYFci5PE](https://www.youtube.com/watch?v=I_XqYFci5PE)

## FP-AI-CXTAWARE1: the best power system saving solution

**MCU**  
**Audio**  
**Scene**  
**Classification**



**Sensor**  
**Human**  
**Activity**  
**Recognition**

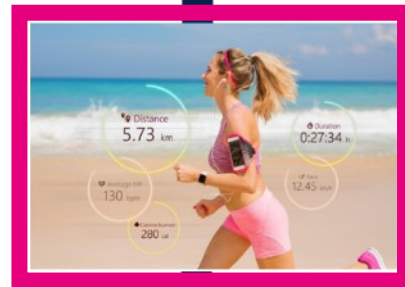
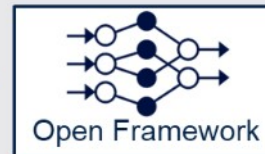
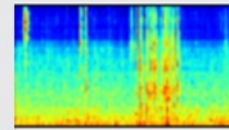


Audio data recording  
Pre-processing  
DB Preparation of NN Training and Test

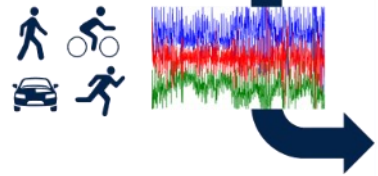
Neural Network Topology  
Definition, Training and Test  
Using existing DL Frameworks

Optimized Neural Network code  
automatically generated for STM32

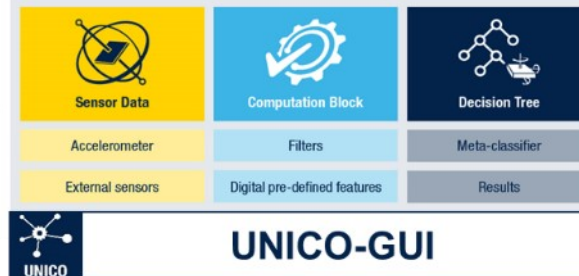
Upload Neural Network  
code on STM32 MCU



Data collection  
Decision Trees generation  
and upload in sensor



Sensor with Machine Learning Core



SensorTile.Box





FP-AI-NANOEDG1



## Step 1 : (PC Side)

### Creation of an ANOMALY DETECTION Machine Learning library

In just a few steps, create a machine learning library\*, custom to your project, and based on a small amount of data captured using your sensor.

PC  
(Win/Linux)

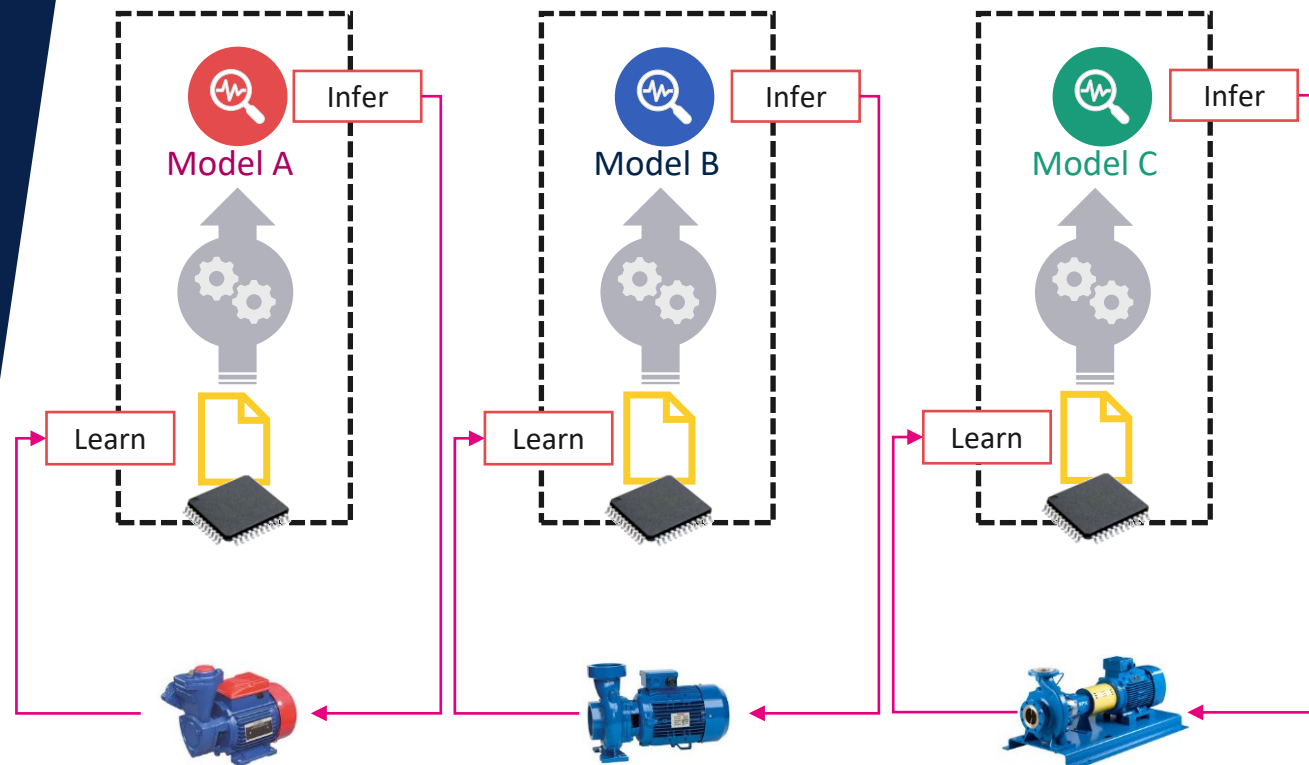


\*A NanoEdge AI Machine Learning Library is a self learning engine, that will train a ML model, inside the Microcontroller and based on locally acquired data/signal.

## Step 2 : (MCU Side)

### Use of an ANOMALY DETECTION Machine Learning library

Each library learns the engine on which it is placed and then analyses it by comparing the learning done locally with the new signals that are coming.







# NanoEdge AI by Cartesiam

NanoEdge AI is a static AI library  
for embedded c software  
running on any Arm Cortex MCU

Learning and inference  
done at the edge.

No pre-trained neural network needed

All work (learning and inference)  
executed inside the STM32 MCU



# SPC5 Studio.AI





# AI beyond autonomous-driving

## Multiple scenarios

### Electrification

- Battery Management: State of Health and Charge
- Hybrid: Efficient Propulsion Mix
- Transmission: Improved Torque Control

### Predictive Maintenance

- Early Detection of anomalies
- Dynamic recalibration to improve robustness

### Vehicle Security

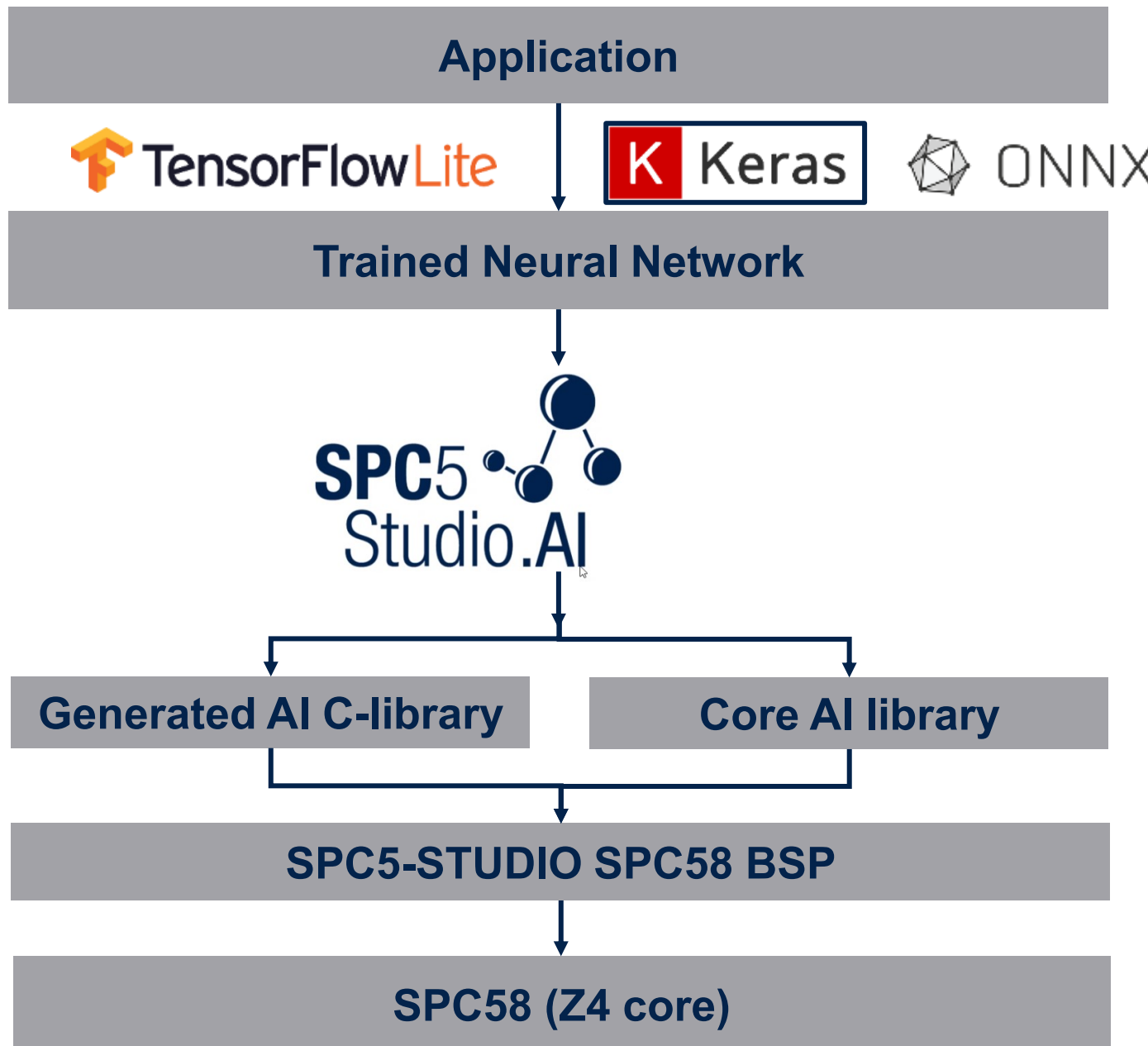
- Hacking, Malicious Attacks
- Physical Vehicle Tampering

### In-Cabin Behaviour Monitoring

- Emergency conditions: Drowsiness, Illness, “Backseat Child”
- Adaptive Comfort settings & driving mode (Sport / Eco)
- Driver identification and monitoring

### Sensor Augmentation

- NOx emission prediction
- Virtual sensors (Vehicle Attitude)
- Sensor degradation compensation



- OpenSource IDE based on ECLIPSE
- Full MISRA 2012 compliant register level access (RLA) low level drivers
- MISRA 2012 checking for customer code
- Visual MCU's pins configuration, run modes and full clock tree configuration with automatic constraints checking
- FreeRTOS support
- Software examples for discovery kits and premium evaluation boards covering most used peripherals
- Compilers: GCC, GHS (green hills software), HighTec EDV-System
- **Artificial Intelligence plug-in; generated software easily portable across different boards**

- 1) **Intrusion detection**
- 2) **Battery charge prediction**
- 3) **Battery lifetime prediction**
- 4) **Audio de-noising**

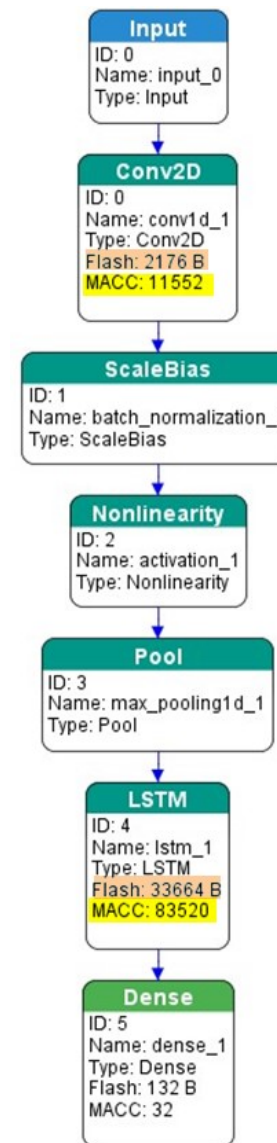


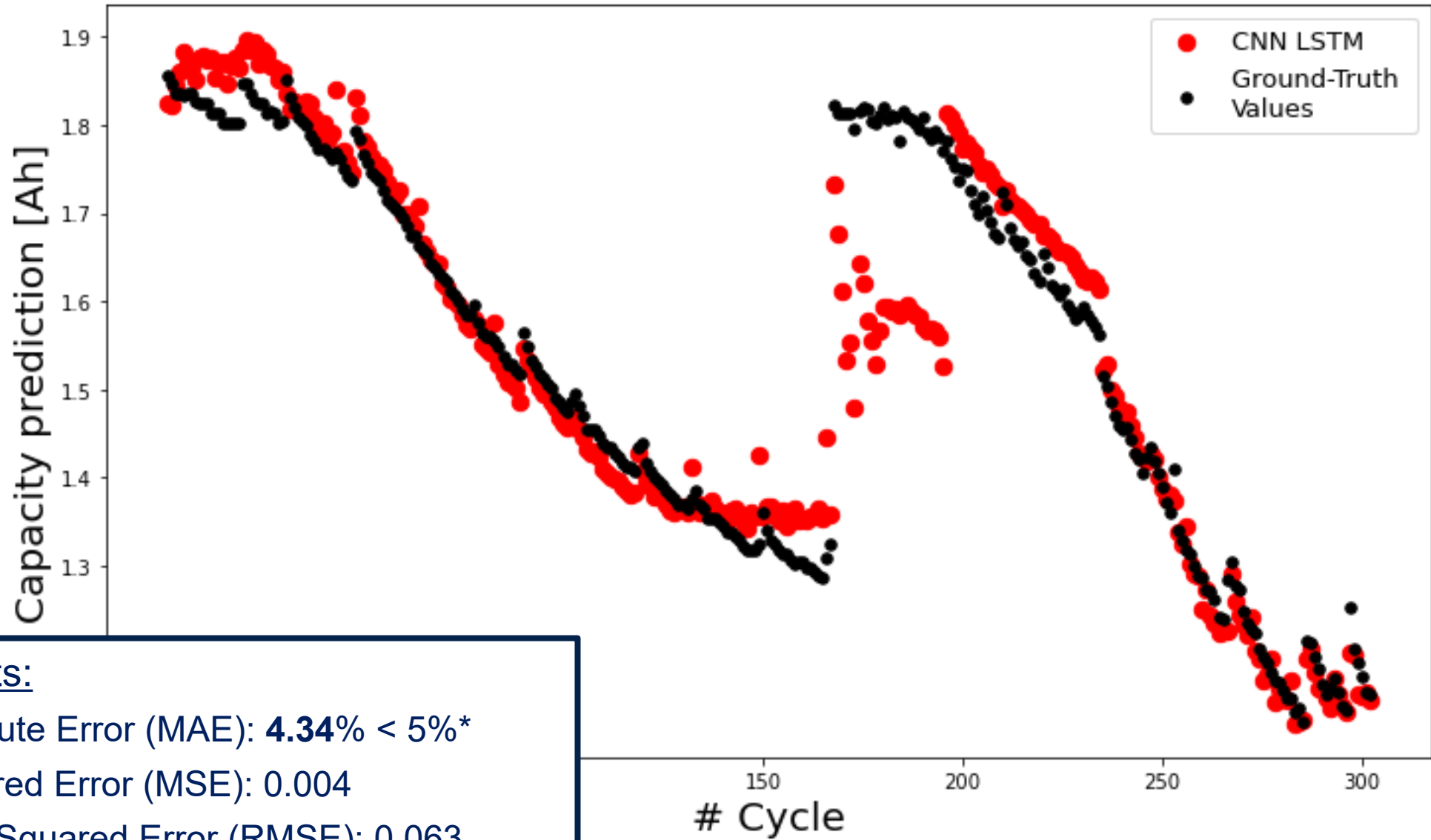
# Battery Management with Long Short-Term Memory (LSTM)

Example Neural Net provided with SPC5 Studio.AI

Actual Performance on Chorus family of automotive MCUs

Measured with SPC5 Studio.AI rel. 2.0.0, 32b-floating-point			
Device	Flash used (kB)	RAM used (kB)	Inference time (ms)
<b>SPC584B</b> 1 (of 1) z4 core @120MHz	149.05	6.66	6.3791
<b>SPC58EC</b> 1 (of 2) z4 @180MHz	152.35	6.79	4.3906
<b>SPC58EG</b> 1 (of 3) z4 @180MHz	154.46	6.91	4.398
<b>SPC58NH</b> 1 (of 3) z4 @200MHz	153.98	6.85	3.7923





### Testing results:

- Mean Absolute Error (MAE): **4.34%** < 5%\*
- Mean Squared Error (MSE): 0.004
- Root Mean Squared Error (RMSE): 0.063

\* EVs acceptable SoH error range

# Design ML



Tiny



=



# Q&A

# Thank you

danilo.pau@st.com

© 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





# Copyright Notice

This presentation in this publication was presented as a tinyML<sup>®</sup> Talks webcast. The content reflects the opinion of the author(s) and their respective companies. The inclusion of presentations in this publication does not constitute an endorsement by tinyML Foundation or the sponsors.

There is no copyright protection claimed by this publication. However, each presentation is the work of the authors and their respective companies and may contain copyrighted material. As such, it is strongly encouraged that any use reflect proper acknowledgement to the appropriate source. Any questions regarding the use of any materials presented should be directed to the author(s) or their companies.

tinyML is a registered trademark of the tinyML Foundation.

[www.tinyML.org](http://www.tinyML.org)