



Fixed AI: Practical Use of Edge AI Unified Core Technology on Sensors and MCU

Danilo PAU

Technical Director

IEEE, ST and AAIA Fellow,

IEEE Distinguished Industrial Lecturer

APSIPA Life, Sigma XI, NAAI member

System Research and Applications

Motivations to hear this talk



The origins of Fixed AI and the dawn of TinyML

The Heterogeneity challenge at the EDGE

Solutions to this challenge

Practical use of the solutions on tiny sensor and MCUs

The cloud returns and Gen EdgeAl





Reducing the Dimensionali ty of Data with Neural Networks G. E. Hinton*

2006

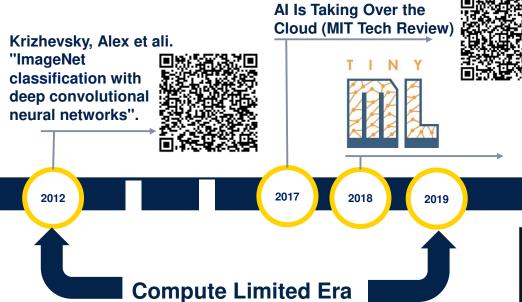
Fixed Function Al

- Low power processing (< mW)
- React on data availability
 - Provide assessments, insights

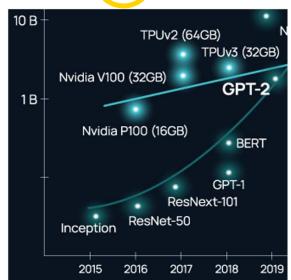
TinyML Products

proliferation

2024



Source https://www.celestial.ai/technology





Tiny Edge

Sensors

Embedded Computing

Actuators



MCU/MPU MCU/MPU

Control









10s KiB RAM No eFLASH KHz to 10 MHz CPU INT1, INT8,INT16; FP32 ≤ 1(o/chip)vs16(off/chip) MiB RAM ≤ 2(o/chip)vs64(off/chip) MiB FLASH ≤ 800 MHz NPU (1 GHz, 4.2 MB o/chip, 600 8b GOPS) INT8,INT16; FP32



≤ 32 KiB RAM ≤ 512 KiB ePCM 160-200 MHz INT8,INT16; FP32









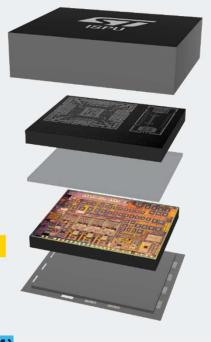
Raw data from 6 DOF MEMS sensor



ISPU processes AI (in a few μW)

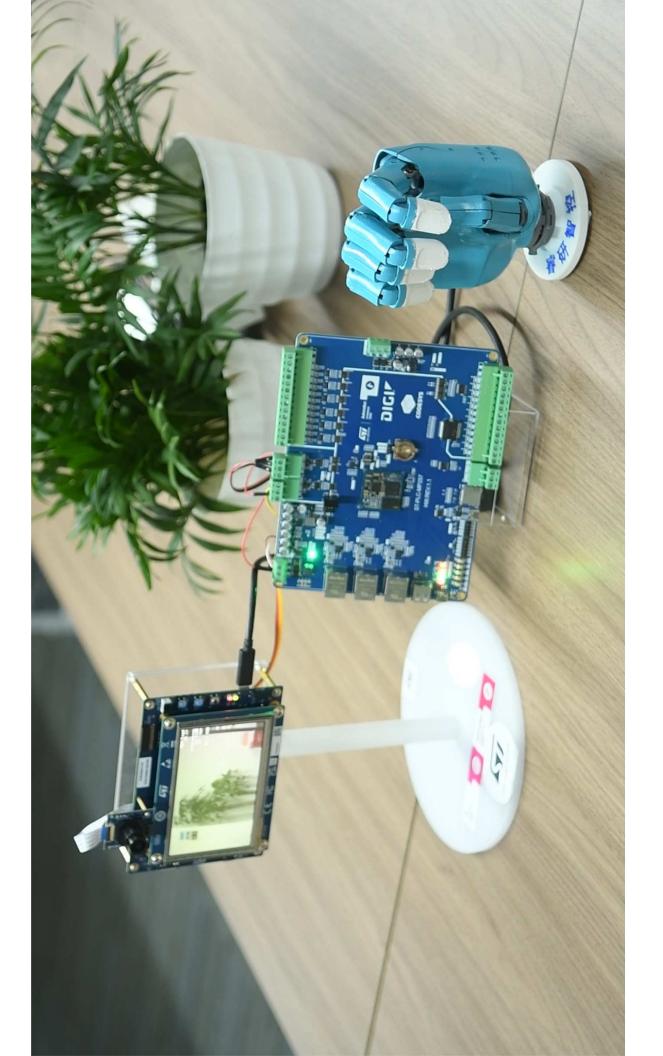


Sensor Hub (up to 4)



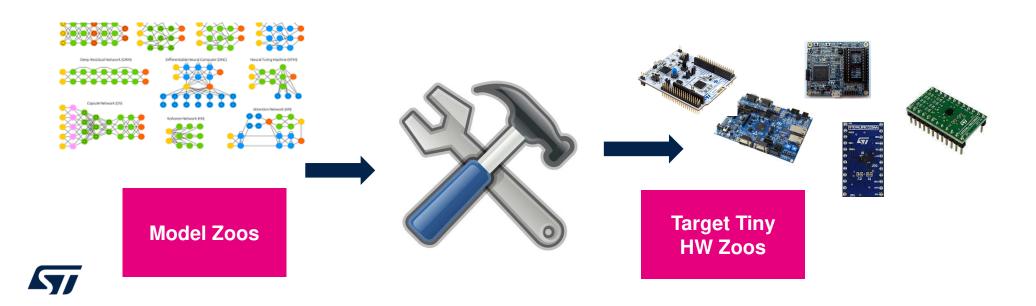




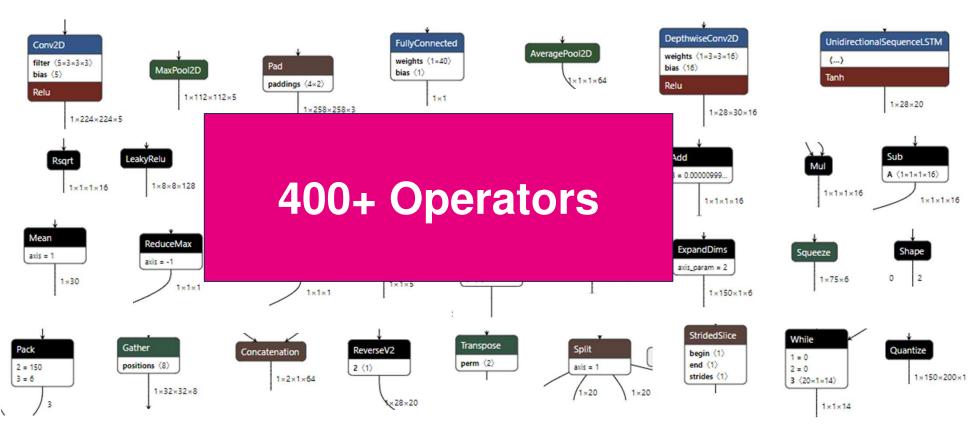


Machine Learning Heterogeneity

- Very different w.r.t.:
 - Source Model: operators, topologies, footprint, representation formats, etc.
 - Execution Target: computational capabilities, available memory, optimized instructions, etc.



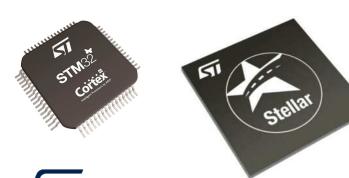
Heterogeneity: Operators





700+ ICs





Heterogeneity: Execution Targets

Instruction Sets (ARM Cortex M 0/4/7/33/55/85, R52, STRed)

Tiny chips, heterogeneous memory schemes, e.g.,

Sensors: MLC, ISPU

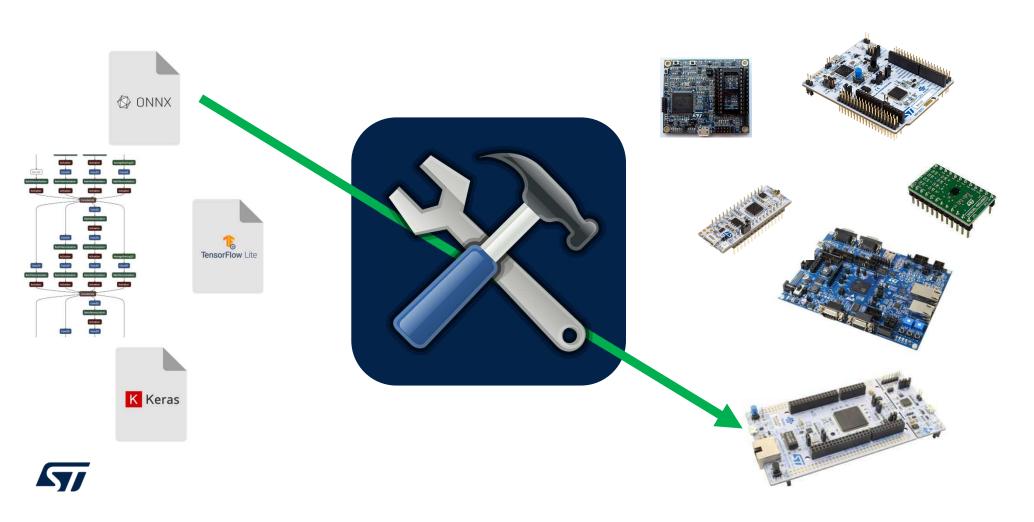
STM32 MCUs, STM32N6, STM32MP2, Stellar MCUs, x86-64

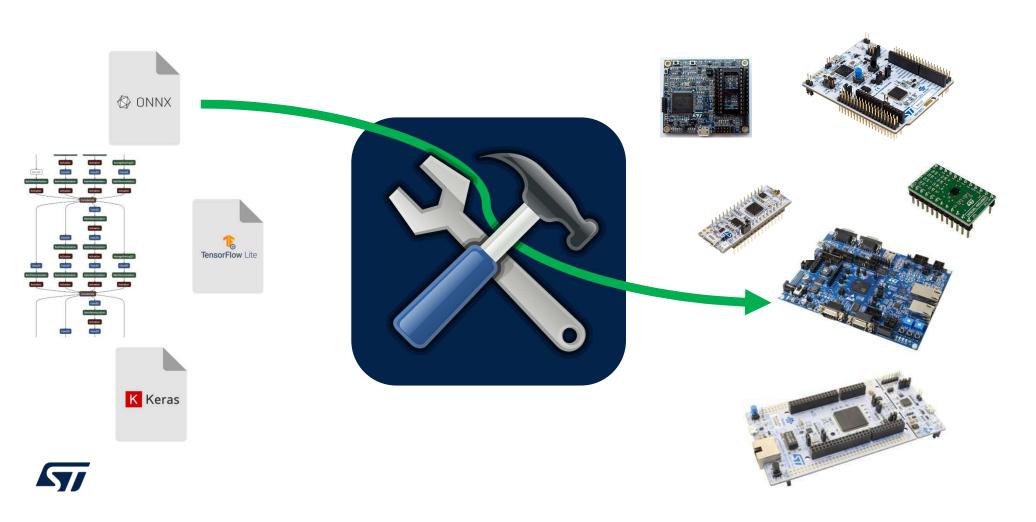
Machine models, e.g.,

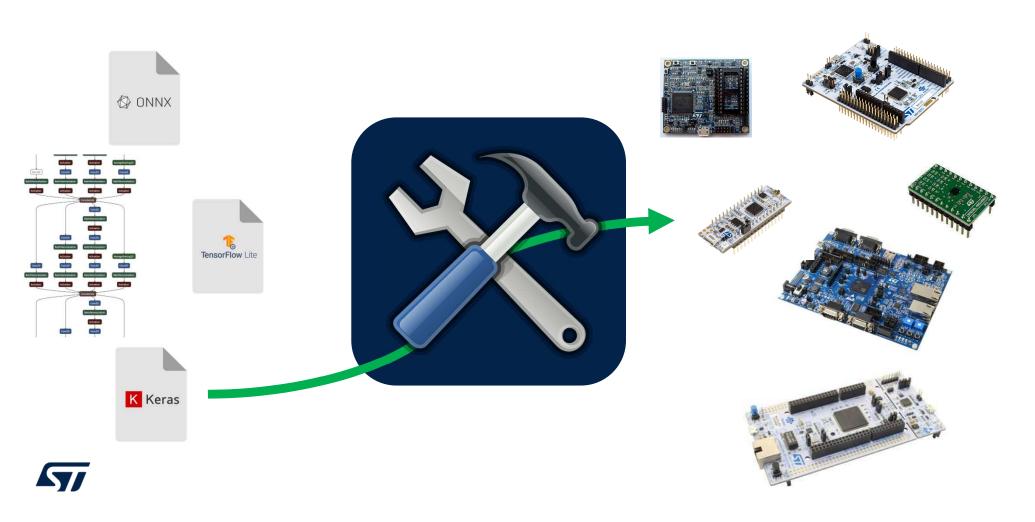
Binary instructions in ISPU Integer SIMD/Vector instructions in ARM Convolutional co-processors and DMAs in STM32N6

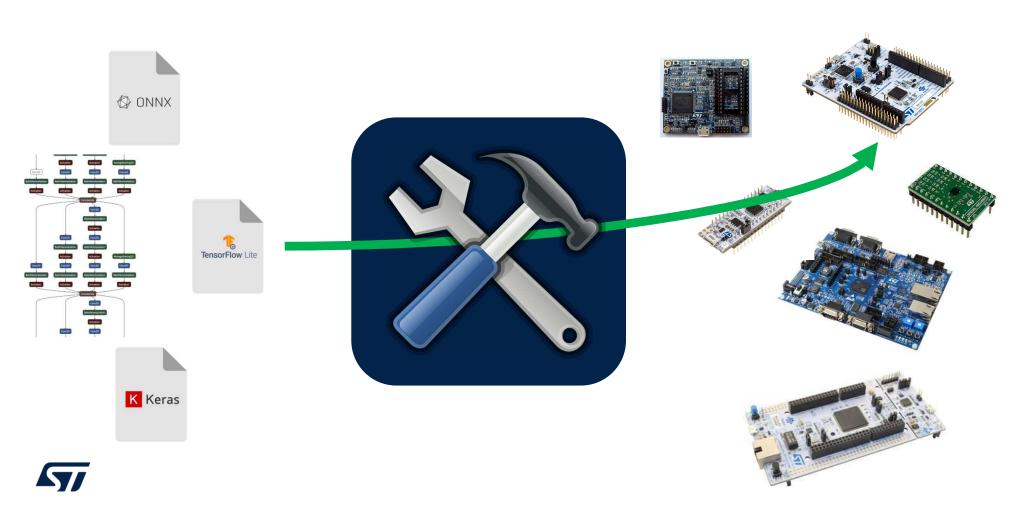
Compilers, e.g.,

gcc, IAR, Hightec, Keil, etc Neural-ART compiler (STM32N6) ISPU compiler

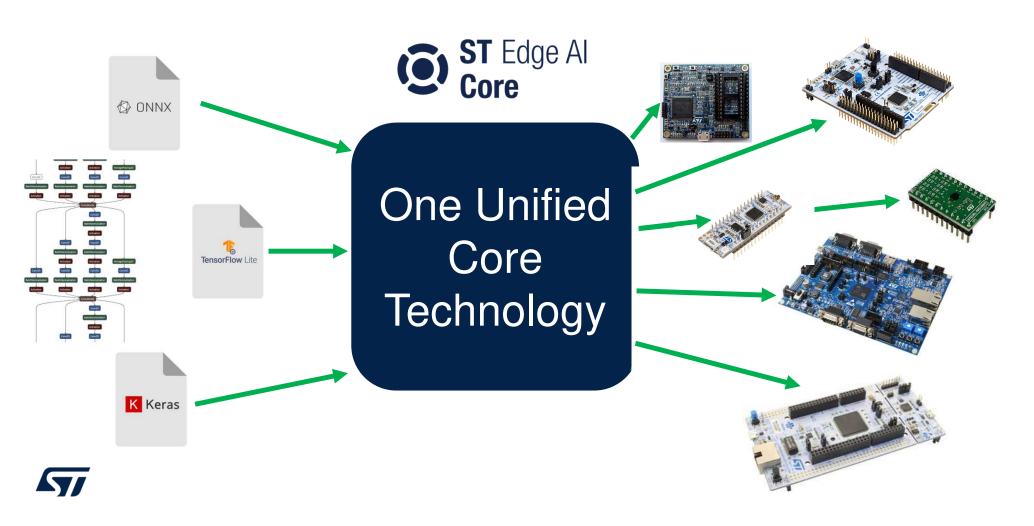






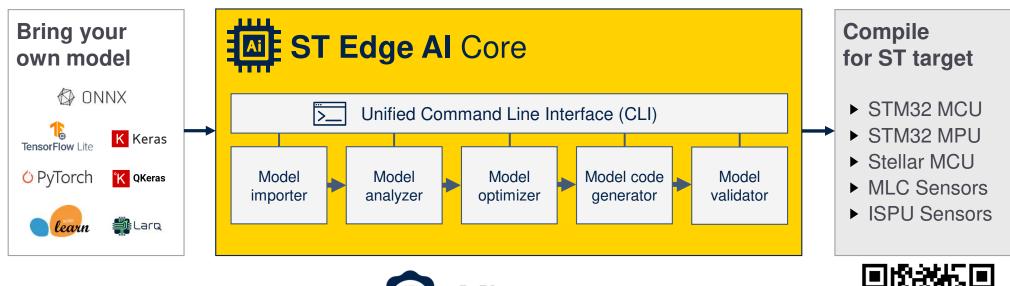


Homogenize Developer Experience



ST Edge AI Core technology

Common state-of-the-art optimizer technology for ALL ST devices

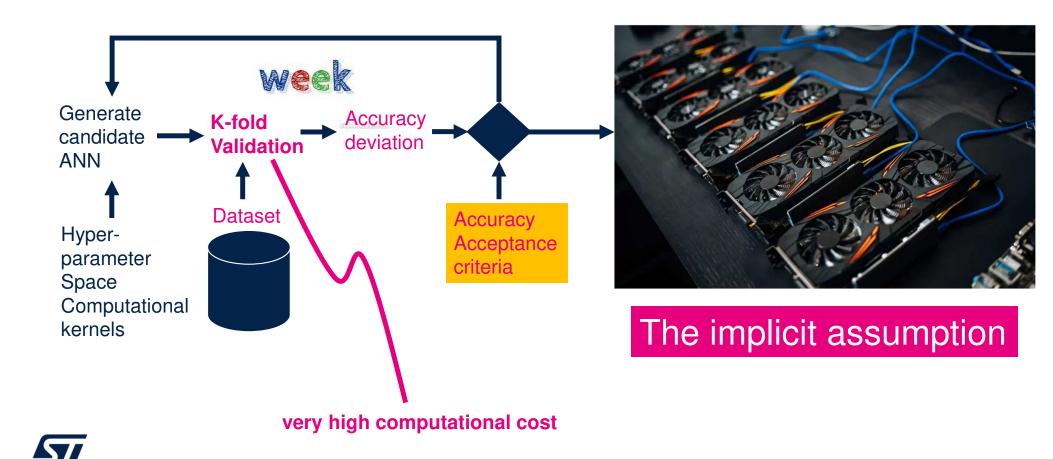






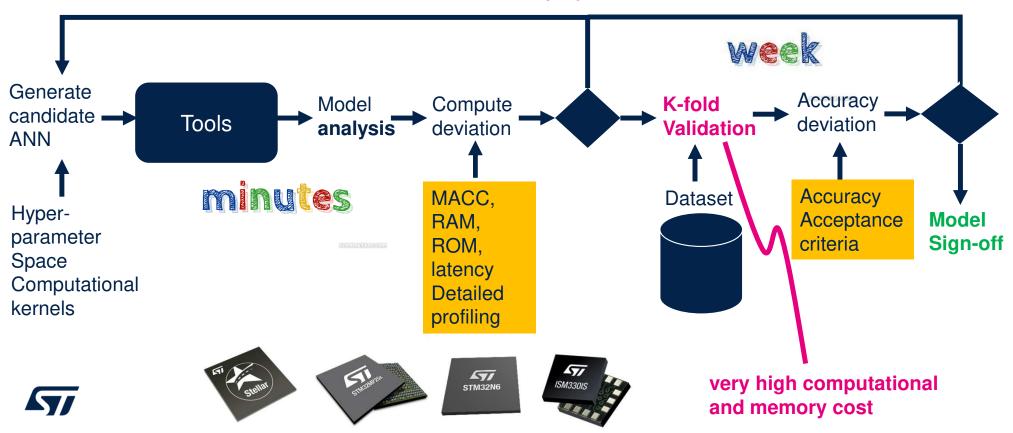


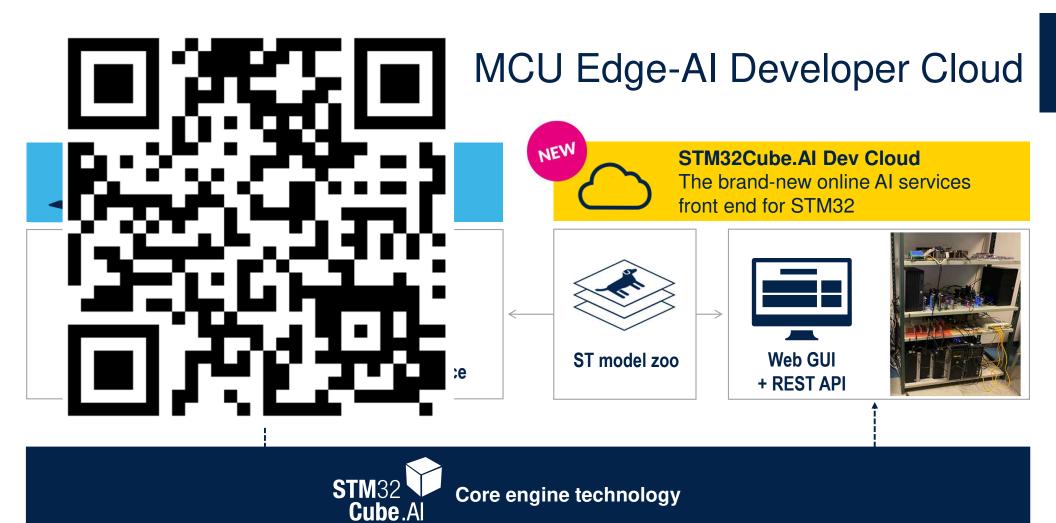
Deployment un-aware NAS/HPO



Deployment aware NAS/HPO

Not deployable on MCU/SENSOR









Fixed Function Al

- Low power processing (< mW)
- React on data availability
- Provide assessments, insights



Attention Is All You Need

- arXiv:1706.03762 Reason, plan, act upon the queries
 - Multi modal
 - Conversational HMI
 - Energy efficient processing (100sT/W)

expanded



Generative Al

danilo.pau@st.com

Our technology starts with You





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.

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

