



Introduzione all'Intelligenza Artificiale nei Microcontrollori

Valeria Tomaselli

ST Confidentia

Tiny Machine Learning Power





Edge AI/TinyML



'A Tsunami of TinyML Devices is Coming' R. El-Ouazzane, STMicroelectronics



3 ST Confidential

Agenda



- 2 Tiny Machine Learning (TinyML) and its benefits
- 3 TinyML use cases and examples

- 5 Training Neural Networks
- 6 TinyML deployment challenges
- 7 ST Edge AI Core STM32Cube.AI

4 Neural Networks

8 Next Steps & Conclusion



Introduction to Artificial Intelligence (AI) & Deep Learning





Al is used today in almost every market segment



life auamente

What is AI?

The evolution of AI





Machine Learning: Why do we need it?

When a **complex task or problem** involves a **large amount of data** and **lots of variables**, but no existing formula or equation can solve it



An example of difficult program

- How to recognize the handwritten digits?
- What makes all these numbers to be identifiable?
- Is there a pattern?
- What is it that makes a 2 to be identified as a 2?



Some examples from MNIST database (Mixed standard institute for standard and technology)



Standard vs Machine Learning algorithm approach



Standard programming "A priori" approach

Design algorithm specific for the given problem



Machine Learning **"Empirical" approach**





Tiny Machine Learning (TinyML) and its benefits





Signals turning into data

Embedded applications will collect more data in the future



Al offers the best approach to process this growing amount of data



Algorithms and predefined models to analyze data and make predictions or decisions



Traditional approaches have limitations:

- when dealing with large datasets
- when the **phenomena are too complex**



Machine learning algorithms to automatically learn patterns and relationships from data





Al-based data processing offers a more flexible and powerful approach to analyzing and making decisions from large data collections

The rise of edge AI



Ultra-low latency Real-time applications

01 Reduced data transmission10 Generate meaningful information

Enhanced privacy and security No data sharing in the cloud 1, 97,



Power efficiency Low-data / Low-power



Improved accuracy analyze data from a wide range of sensors and sources

Edge AI benefits many application domains:

Industrial maintenance Condition monitoring Predictive maintenance



Control systems From home heating systems to industrial machines



Internet of Things (IoT) Smart cities, smart buildings, connected homes, and industrial automation



Edge AI adoption is rapidly accelerating



life.augmented

Source: ST customer survey – December 2023 14 ST Confidential

TinyML use cases and examples







ISPU for personal application

Al at the edge with ultralow power 6-axis IMU for consumer market





A completely new level of capabilities and detection accuracy in human activity recognition applications:

- Consumer health •
- Gesture recognition •
- Activity recognition ٠
- Motion tracking ٠

Carry

position



Gait

analysis



Pose

estimation



Fall

detection







Active time



Fitness

activities





Activity

recognition







recognition

And more...





ISPU to approach the Industry 5.0

Al at the edge with ultralow power 6-axis IMU for industrial market





Higher detection accuracy, always on monitoring in anomaly detection applications

Home alarms

The second

- Robotics
- Condition monitoring









AI Solutions on STM32

A full development ecosystem to create your AI application



Al extension for STM32CubeMX to map pre-trained Neural Networks



Person presence detection Food classification



STM32 **Community** with dedicated Neural Networks topic and **AI expert partners**



People activity recognition Audio scene classification





life.gugmente

Trainings, hands on, MOOCs and partners **videos**



Condition-based monitoring





Automotive Edge AI Solutions

Al at the edge with Stellar E microcontroller





Better efficiency & maintenance

- Protecting critical functions
- Improving energy management
- Enabling predictive maintenance

A personalized driving experience

- optimize comfort and convenience
- power infotainment systems





ST Confidential





What are Neural Networks?

- Also referred to as Artificial Neural Networks.
- Inspired by human neural system.
- Human neuron has three main components
 - Dendrites
 - Take inputs from other neurons in terms of electrical pulses.
 - Cell body
 - Makes the inferences and decides the actions to take.
 - Axon terminals
 - Send the outputs to other neurons in terms of electrical pulses.







ST Confidential



• The heart of a neural network



23





²⁵ ST Confidential



²⁶ ST Confidential





Neural Networks

Notations





Different types of Neural Networks

Multi Layer Perceptrons

Convolutional Neural Networks

Recurrent Neural Networks









Convolutional Neural Network







Confidential

2D Convolutions





3D Convolutions

Input Volume (+pad 1) (7x7x3) Filter W0 (3x3x3)													
x[:,:,0] w0[:,:,0]													
0	0	0	0	0	0	0	-1 0 1						
0	0	0	1	0	2	0	0 0 1						
0	1	0	2	0	1	0	1 -1 1						
0	1	0	2	2	0	0	w0[:,:,1]						
0	2	0	0	2	0	0	-1 0 1						
0	2	1	2	2	0	0	1 -1 1						
0	0	0	0	0	0	0	0 1 0						
x[:::1] w0[:::,2]													
0	0	0	0	0	0	0							
0	2	1	2	1	1	ø	1 1 0						
0	2	1	2	0	1	ø	0 -1 0						
0	0	2	1	0	1	0	Bias b O(1x1x1)						
0	1	2	2	2	2	0/	b0[:,:,0]						
0	0	1/	2	0	1/	6	1						
0	V	0	0	0/	6	0							
xí:		21	7	/		1							
0	0	0	0	0	0	0							
0	2/	1	1	2⁄	0	0							
9	1	0	ø	1	0	0							
0	0	1	0	0	0	0							
0	1	0	2	1	0	0							
0	2	2	1	1	1	0							
0	0	0	0	0	0	0							
_													

Filter W	'1 (3x3x3)	Output Volume (3x3x2)						
w1[:,	:,0]	0[:,:,0]						
0 1	-1	2	3	3				
0 -1	0	3	7	3				
0 -1	1	8	10	-3				
w1[:,:,1] o[:,:,1]								
-1 0	0	-8	-8	-3				
1 -1	0	-3	1	0				
1 -1	0	-3	-8	-5				
w1[:,	:,2]							
-1 1	-1							
0 -1	-1							
1 0	0							
Bias b1 b1[:,	(1x1x1) :,0]							

toggle movement







- Unit step
 - Threshold
- Sigmoid Function
 - Like a step function but smoother
 - Best to predict probabilities
- Tan hyperbolic
 - Stretched out version of the sigmoid function
- ReLU
 - Computationally efficient





- Function choice depends on the characteristics of the data.
- For example Sigmoid Function works good for classification purposes, resulting in Faster training and convergence.
- ReLU is good for approximation.
 As it is simple so always start from this if you don't know the data characteristics. Helps against gradient vanishing
- We can also define custom activations.



Learning hierarchical representations

• More than one stage of non-linear feature transformation





Feature visualization on convolutional net on ImageNet [Zeiler & Fergus 2013]

Recurrent Neural Networks

- Convolutional neural networks have no internal state persistence
- Recurrent neural networks address this issue. They are networks with loops in them, allowing information to persist.



Unrolling recurrent neural network loops


Training Neural Networks

, (ч) О 0 iuu xl



Training neural networks

- In supervised learning an assumption is to have a relatively large labeled dataset.
- Feed all the samples as inputs to get an output. Called forward propagation or inference run outputs.
- At start the weights can be randomized or predefined depending on the applications scenario.
- The result \hat{y} is compared with ground truth output y.
- The task is to make the output value ŷ to be as close to y as possible reducing the error expressed as Loss functions L(ỹ, y).





Training neural networks





life.augmented



- Go back and adjust the weights slowly. Aim is $Error_T < Error_{T-1}$
- Repeat this process until the error we get is very small.

 $\lim_{\epsilon \to 0} \ Error_T < \in$











- Brute force
 - Try all the possible combination of weights.
 - Plot the cost function.
 - Use the weights which result in smallest error.
 - Sounds simple but will take too much time !!!

41

• Enters the gradient descent

Gradient Descent



• Enters the gradient descent

Gradient Descent





Backpropagation pitfall









Note: One epoch means one pass of the whole training set.

Gradient descent in action

• Example: Finding best linear fit to a set of points.





Under, Good, and Over Fit



life.augmented

47

Learning datasets

- One of the challenges of working with ANN is to have a big labelled dataset.
- The dataset is usually divided into
 - Training set
 - Training set
 - Validation set
 - Testing set





Learning Curves



Example of Training Learning Curve Showing an Underfit Model That Requires Further Training

Example of Train and Validation Learning Curves Showing an Overfit Model

Example of Train and Validation Learning Curves Showing a Good Fit



TinyML deployment challenges





Cloud computing vs Edge Computing

Cloud Computing



Training is still here

Edge Computing





And many more...
Inference can be here



What are the main deployment challenges of TinyML?







Tiny ML Challenges



Phenomenal Cosmic Powers!

Itty Bitty Living Space!

ST Confidential

53

Tiny ML Challenges



Tiny Resources

Toward Zero Power

Achieve High Accuracy

Live without floating point numbers

Automated deployments



5 Key steps for Supervised Deep Learning



Deploying Inference on MCU/sensor





"Embedded Real-Time Fall Detection with Deep Learning on Wearable Devices," 2018 21st Euromicro Conference on Digital System Design (DSD), Prague, 2018, pp. 405-412, doi: 10.1109/DSD.2018.00075.

Deployment un-aware Neural Architecture Search (NAS)





very high computational cost

Deployment aware Neural Architecture Search (NAS)

Not deployable on MCU/SENSOR



ST Confidential

Solving Heterogeneity





Machine Learning Heterogeneity





Coping with source and target heterogeneity































ST Confidential

















1

ST Confidential





1

ST Confidential
Heterogeneity: Operators





Source Model Heterogeneity: DL Formats



.tflite























Availability of operators: not all operators are available in all frameworks

Attributes: different attributes, meaning and default values

Data layout, i.e., channel first vs. channel last

Quantization, i.e., many possible bit-width to represent data





TFLite MEAN ONNX Mean











































Limiting to binary, int8, float32 there are 3^3=27 combinations to be supported for a single layer

Output



Input

Heterogeneity: Execution Targets

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

Different computational power, e.g., MLC, ISPU STM32 std MCU family, STM32N6

Different specialized functional units, e.g., Binary accelerator in ISPU Integer SIMD/Vector ISA in ARM Convolution accelerators in STM32N6



Heterogeneity: Execution Targets



Memory optimizer

Optimize memory allocation to get the best performance while respecting the constraints of the embedded design

- Memory allocation
- Internal/external memory repartition
- Model-only update option

- Different memory infrastructure e.g.,
 - Single memory component
 - Multiple and homogeneous memories
 - Multiple and heterogeneous memories





life.augmented





life.augmented



life.augmented

Homogenize Developer Experience



life.augmente

Removing Deep Learning Frameworks Differences



Removing Deep Learning Frameworks Differences



Removing Deep Learning Frameworks Differences





Hiding Target Differences



- Same public APIs for all the execution targets
 - Same application code means portability on all the supported target
- Internal code is optimized for the different targets exploiting the heterogeneous hardware capabilities
- Extra APIs to support advanced features of hardware



Hiding Target Differences



- Simplified public APIs: Init + run + deinit
- Different Optimization
 Objectives
 - **Time**: to minimize inference time
 - RAM: to minimize use of memory
 - **Balanced**: trade-off between inference time and memory usage
- Multiple networks instancing in the same application



ST Edge Al Core





More than a decade of research, development, and deployment





*Intelligent Sensor Processing Unit
Announcing the ST Edge AI Suite

The most complete developer-centric approach to accelerate the deployment of edge AI





ST Edge AI Core



life.augmented

STM32Cube.Al Edge Al optimization tool for STM32

	Optimus Colorada Benchmark Results Optimus Colorada Colorada Colorada Colorada Colorada Colorada Colorada	
	Model currently selected	
X	Nepur Output MODELTYPE MACC Sold 1-224-024-0 Meet D-1-1-18 Meet 1227/1984	A loss of the
	Select year model optimization options Select	
	Since of epidemization results (other ness time bits well be measured in bandwall bits performance) Since Terminal C C Dres - 0 Operation Operat	



STM32Cube.AI solutions



life.augmented

113 ST Confidential

STM32Cube.AI: your go-to solution for edge AI

Fully **integrated** in the STM32Cube software development suite

Cutting-edge **performance**, optimized for STM32 hardware



Completely **FREE** of charge



STM32Cube.AI core technology

Optimize and validate your NN model





Two versions of the same tool depending on your profile





STM32 benchmarking tool

The unique possibility to evaluate the performance of models remotely, on real STM32 boards



Get the real inference time from optimized models running on STM32



Benchmark models on a large variety of STM32 boards

Find the most appropriate board for your application



Get access to the most recent devices A board farm is constantly updated with the latest available boards





STM32 model zoo

A collection of application-oriented models optimized for STM32





Hosted on GitHub



Model training scripts

Scripts to generate and validate



Getting started application packages

- Automatically generated from the trained models
- Easy to deploy for end-to-end evaluation



The 3 pillars of STM32Cube.Al

Graph optimizer

Automatically improve performance through graph simplifications & optimizations that benefit STM32 target HW architectures



- Auto graph rewrite
- Node/operator fusion
- Layout optimization
- Constant-folding...
- Operator-level info to finetune
 memory footprint and computation

Quantized model support

Import your quantized ANN to be compatible with STM32 embedded architectures while keeping their performance



- From FP32 to Int8 or mixed-precision
- Minimum loss of accuracy
- Code validation on target
 - Latency
 - Accuracy
 - o Memory footprint

Memory optimizer

Optimize memory allocation to get the best performance while respecting the constraints of your embedded design



- Memory allocation
- Internal/external memory repartition
- Model-only update option

STM32Cube.AI is free of charge, available both in graphical interface and in command line.



Graph optimizer

Squeeze your graph to fit into an MCU!





Fully automated process in the STM32Cube.Al workflow

- Your original graph is optimized at the very early stage for optimal integration into the STM32 MCU/MPU
- Loss-less conversion



Quantized model support

Simply use quantized networks to reduce memory footprint and inference time



STM32Cube.AI supports quantized neural network models with **all parameter formats**:

- FP32
- Int8
- Mixed binary Int1 to Int8 (Qkeras*, Larq.dev*)

*Please contact <u>edge.ai@st.com</u> to request the relevant version of STM32Cube.AI



HW Target: NUCLEO-STM32H743ZI2 Model: Low complexity handwritten digit reading Freq: 480 MHz Accuracy: >97% for all quantized models

Tested database: MNIST dataset





Memory optimizer

Optimize performance easily with the memory allocation tool

	network			
			conv2d_3	
HGB			id: Type: MACC: Input Tensor	3 conv2d 82960
SKGB 1 1 1 1 1 1 1 1 1			conv2d_2_output Size: Format: Shape:	36864B int/us (48, 48, 16
KIB L L L			Output Tensor	
KGB			com/2d_3_output Size: Format: Shape:	9216B int/us (24, 24, 16
ві		i i i i i	Scratch Tensor	
0 1 2	4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 2	25 26 27	conv2d_3_scratch0 Size	436B

Model RAM consumption per layer

 Easily identify the most critical layers

	/ ☑ Use external flash Memory: Custom ✓
Model memory allocation	Split weights between internal and external flash using a linker script \checkmark
 Set your external memory 	Start Address: 0x00000000 Size (Mbytes)
 Map in non-contiguous internal flash section Partition internal vs external flash 	Tensor Size Internal 440KB External 0KB conv1_weights 864 ✓ □ conv_bias 32 ✓ □ conv_dw_1_weights 288 ✓ □ conv_dw_1_bias 32 ✓ □ conv_dw_1_bias 32 ✓ □ conv_pw 1_weights 512 ✓ □
memories	Use external RAM Memory: Custom Start Address: 0x0000000
Re-use model input buffer to store activation data*	✓ Use activation buffer Start Address: 0x00000000 Act. size (by 752712 □ Copy weight to RAM Start Address: Weight size: 451496
 Minimize RAM requirements 	✓ Use activation buffer for input buffer (allocate-inputs) □ Force classifier validation output (classifier) ✓ Use activation buffer for the output buffer (allocate-outputs)
Relocatable network	 Split weights during code generation (split-weights) Generate relocatable network (relocatable)
 A separate binary is generated for the library and the network to enable standalone model upgrade 	Report's ouput directory C:\Users\richardv\.stm32cubemx Browse
, , , , , , , , , , , , , , , , , , ,	Enable custom layer support Custom Layer JSON File: Browse Browse



We provide everything to kick off your project

Design documentation





Development zone Get started on application development and project sharing

- Wiki by ST is a great forum to learn and start deploying edge AI on STM32!
- Videos of application examples
- Massive open online courses (MOOCs)

Hardware and software tools



- Evaluation platforms for STM32 MCUs and MPUs
- Additional sensor boards
- Full software suite

Support and updates



- ST Community: STM32 ML & AI group
- Distributor certified FAE
- Support center
- Newsletter



Next Challenges



On-Device Learning

Zero Code

Generative AI on the Edge



Question Session





System Development with STMicroelectronics 125 ST Confidential

Conclusion





ST Confidential

Our technology starts with You



© 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 <u>www.st.com/trademarks</u>. All other product or service names are the property of their respective owners.

