

tinyML[®] Talks

Enabling Ultra-low Power Machine Learning at the Edge

“Suitability of Forward-Forward and PEPITA Learning to MLCommons-Tiny benchmarks”

Danilo Pau – Technical Director, IEEE, AAIA & ST Fellow STMicroelectronics

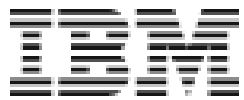
September 19, 2023



www.tinyML.org



Thank you, **tinyML Strategic Partners**,
for committing to take tinyML to the next Level, together



T I N Y



TALKS
webcast

Executive Strategic Partners

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/IIoT



Automotive



Mobile



Accelerate Your Edge Compute

SYNTIANT

Making Edge AI A Reality

www.syntiant.com

Platinum Strategic Partners



**DEPLOY VISION AI
AT THE EDGE AT SCALE**

SONY

Gold Strategic Partners



AHEAD OF WHAT'S POSSIBLE™



AHEAD OF WHAT'S POSSIBLE™

Where what if
becomes what is.

Witness potential made possible at analog.com.

Build the
Future of tinyML

on **arm**



T I N Y



TALKS
webcast



EDGE IMPULSE

The Leading Development Platform for Edge ML

edgeimpulse.com

Decarbonization

Digitalization



Driving decarbonization and digitalization. Together.

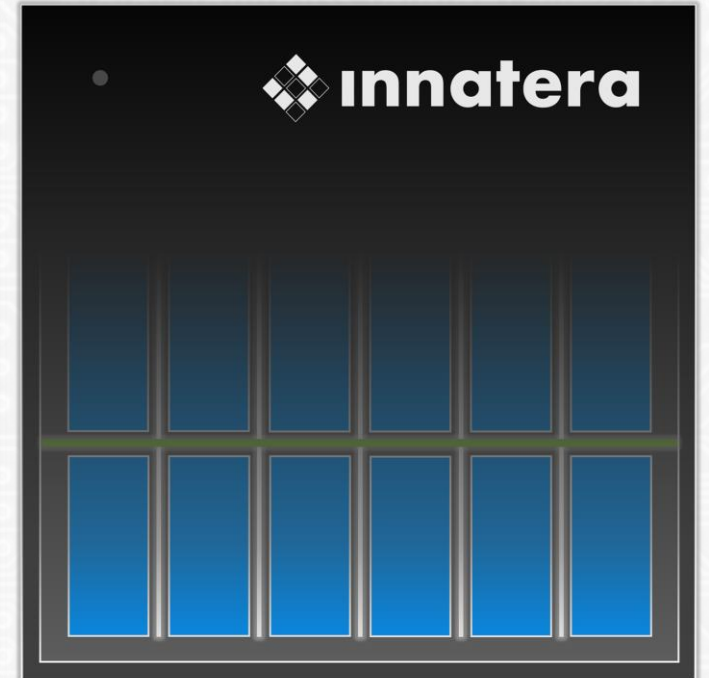
Infineon serving all target markets as
Leader in Power Systems and IoT

www.infineon.com





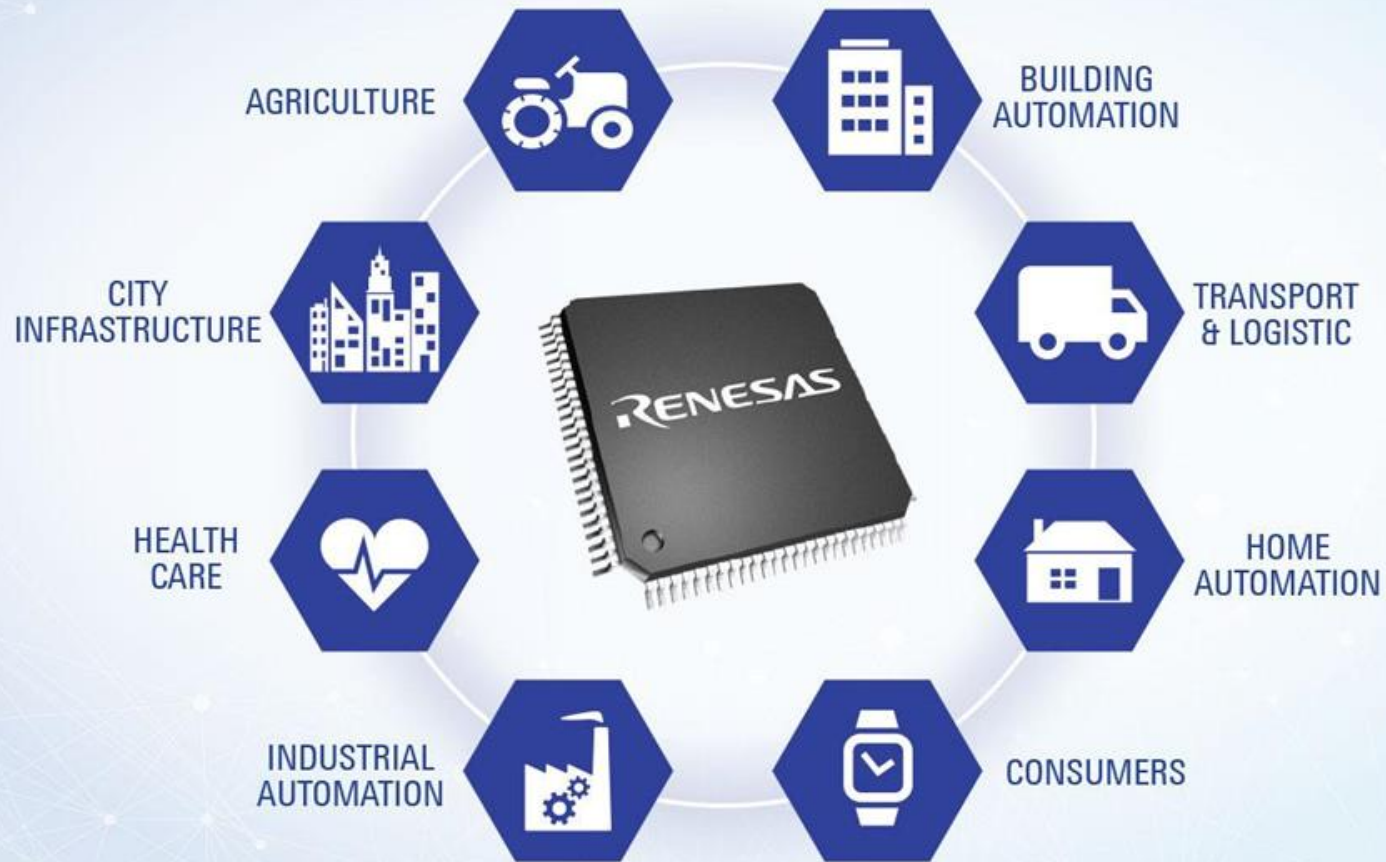
NEUROMORPHIC INTELLIGENCE FOR THE SENSOR-EDGE





Microsoft

Renesas is enabling the next generation of AI-powered solutions that will revolutionize every industry sector.



[renesas.com](https://www.renesas.com)



life.augmented

STMicroelectronics provides extensive solutions to make tiny Machine Learning easy



ENGINEERING EXCEPTIONAL EXPERIENCES

We engineer exceptional experiences for consumers in the home, at work, in the car, or on the go.

www.synaptics.com



T I N Y



Silver Strategic Partners



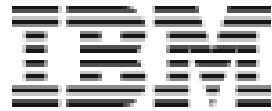
brainchip



GREENWAVES
TECHNOLOGIES



⚡ Grovety Inc.



NotaAI





Join Growing tinyML Communities:



16.8k members in
49 Groups in 41 Countries

tinyML - Enabling ultra-low Power ML at the Edge

<https://www.meetup.com/tinyML-Enabling-ultra-low-Power-ML-at-the-Edge/>



4k members
&
13k followers

The tinyML Community

<https://www.linkedin.com/groups/13694488/>





Subscribe to
tinyML YouTube Channel
 for updates and notifications
(including this video)

www.youtube.com/tinyML



tinyML
4.33K subscribers

10.4k subscribers, 627 videos with 378k views

HOME VIDEOS PLAYLISTS COMMUNITY CHANNELS ABOUT

106 views · 4 days ago	138 views · 4 days ago	54 views · 4 days ago	47 views · 4 days ago	132 views · 4 days ago	137 views · 4 days ago
122 views · 4 days ago	262 views · 2 weeks ago	511 views · 3 weeks ago	229 views · 3 weeks ago	265 views · 3 weeks ago	286 views · 1 month ago
351 views · 1 month ago	462 views · 2 months ago	374 views · 2 months ago	133 views · 2 months ago	287 views · 2 months ago	336 views · 2 months ago
378 views · 2 months ago	214 views · 2 months ago	448 views · 2 months ago	159 views · 2 months ago	190 views · 2 months ago	545 views · 2 months ago



tinyML Asia Technical Forum

**November 16, 2023
Seoul, South Korea**



Register now
<https://www.tinyml.org/event/asia-2023/>

2023 Edge AI Technology Report

The guide to understanding the state of the art in hardware & software in Edge AI.



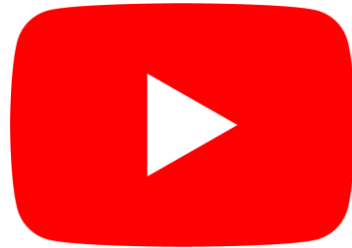


Reminders

Slides & Videos will be posted tomorrow



tinyml.org/forums



youtube.com/tinyml



Please use the Q&A window for your questions





Danilo Pau



Danilo PAU (h-index 26, i10-index 69) graduated in 1992 at Politecnico di Milano, Italy. One year before his graduation, he joined SGS-THOMSONS (now STMicroelectronics) as interns on Advanced Multimedia Architectures, and he worked on memory reduced HDMAC HW design. Then MPEG2 video memory reduction. Next, on video coding, transcoding, embedded 2/3D graphics, and computer vision. Currently, his work focuses on developing solutions for tiny machine learning tools.

Since 2019 Danilo is an IEEE Fellow and AAIA on 2023; he served as Industry Ambassador coordinator for IEEE Region 8 South Europe, was vice-chairman of the “Intelligent Cyber-Physical Systems” Task Force within IEEE CIS, was IEEE R8 Afl member in charge of internship initiative. Today he is a Member of the Machine Learning, Deep Learning and AI in the CE (MDA) Technical Stream Committee CESoc. He was AE of IEEE TNNLS. He wrote the IEEE Milestone on Multiple Silicon Technologies on a chip, 1985 which was ratified by IEEE BoD in 2021 and IEEE Milestone on MPEG Multimedia Integrated Circuits, 1984-1993 which was ratified in 2022. He served as TPC member to TinyML EMEA forum and is the chair of the TinyML on Device Learning working group. He serves as 2023 IEEE Computer Society Fellow Evaluating Committee Members

With 78 and 68 respectively European and US application patents, 157 publications, 113 ISO/IEC/MPEG authored documents and 67 invited talks/seminars at various Universities and Conferences, Danilo's favorite activity remains supervising undergraduate students, MSc engineers and PhDs.



Fabrizio Aymone



Fabrizio M. Aymone is currently pursuing a Bachelor degree in Electronics Engineering at Politecnico di Milano. He is also intern at the System Research and Applications department of STMicroelectronics, where he is studying solutions for On-Device Learning in the domain of tiny devices. His research interests focus on reducing memory usage and computational complexity of AI algorithms and exploring alternative learning rules to backpropagation.



life.augmented

Suitability of Forward-Forward and PEPITA Learning to MLCommons-Tiny benchmarks

Danilo P. Pau and Fabrizio M. Aymone

Sept, 19 2023

Agenda

1 Back-propagation

2 On-Device Learning

3 Forward-Forward and PEPITA

4 The research question

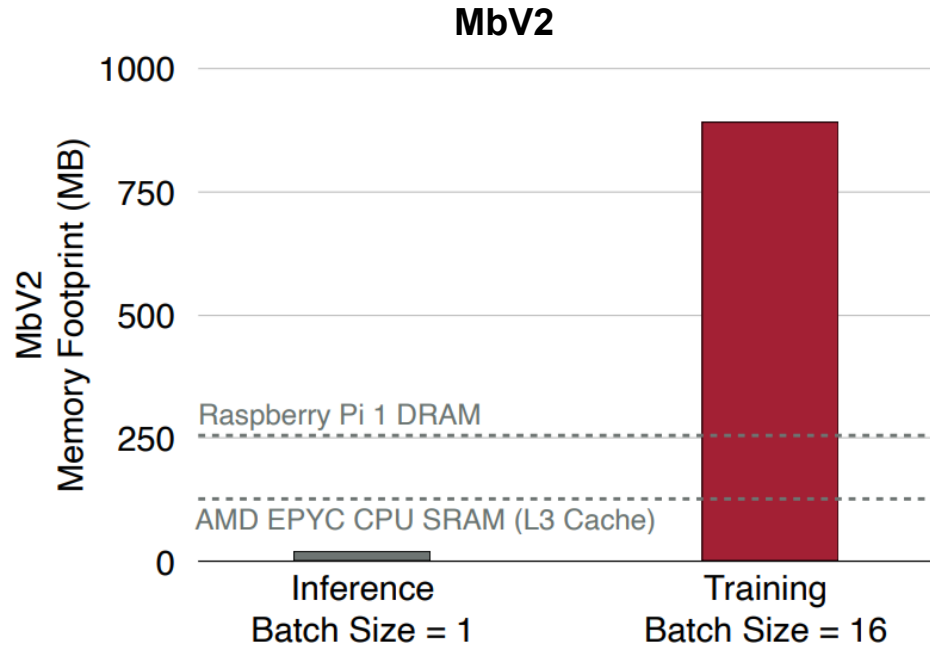
5 Methodology and Results

6 Takeways

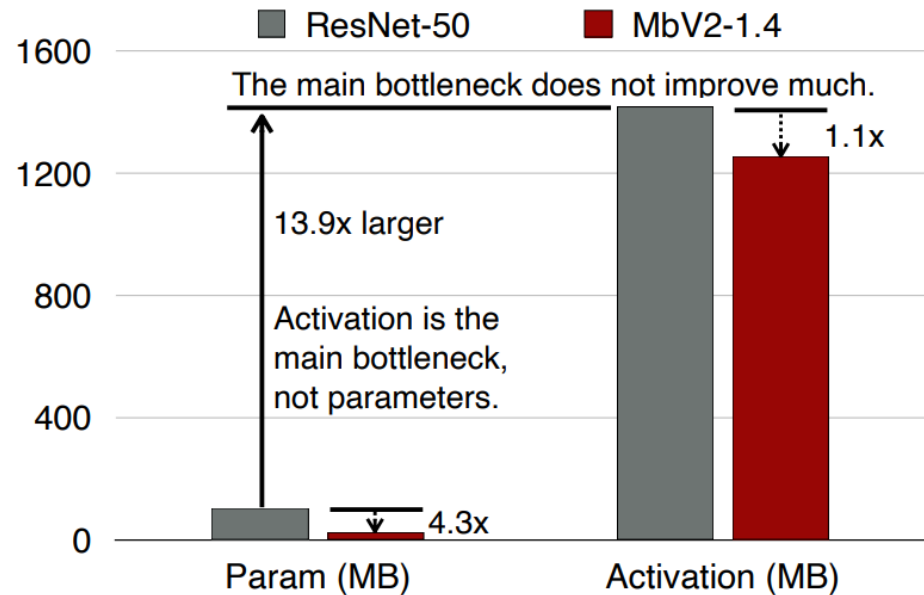
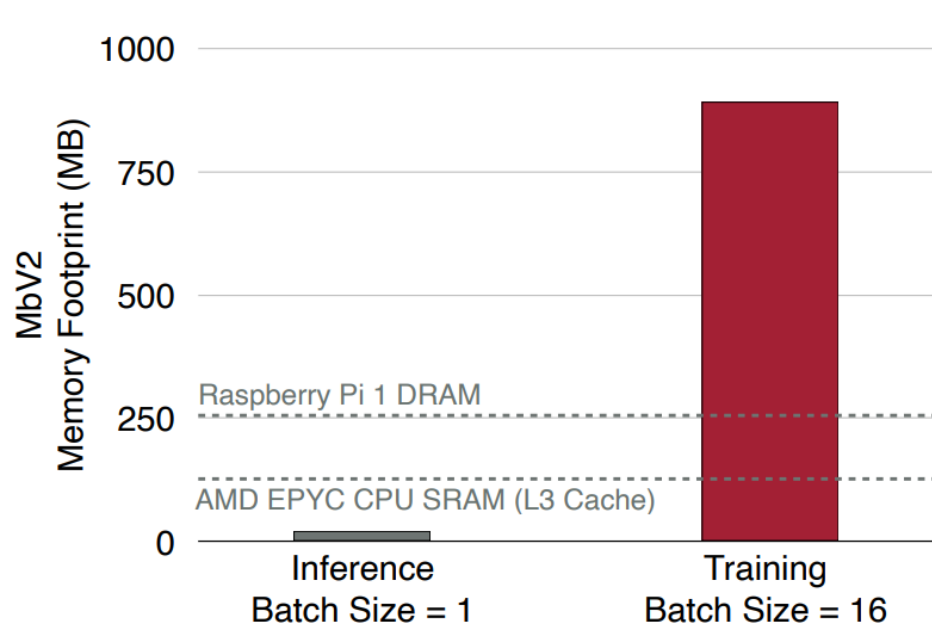
7 Future works

8 Q&A

Reduce activations, not trainable parameters for efficient on-device learning¹



Reduce activations, not trainable parameters for efficient on-device learning¹



Backpropagation

Activations are computed for each layer and stored into memory

- The loss is calculated w.r.t. the ground truth and the final output;
- The gradient of the loss is then computed;
- Then, the derivative of the output activations z_l
- Finally, the derivative of the loss function for the output layer is computed as Hadamard product

Algorithm 1 Backpropagation

Forward Pass

for $l = 1, \dots, L$ do

$$a_l = \sigma_l(W_l a_{l-1} + b_l)$$

end for

Backward pass layer L

$$\delta_L = \nabla_{a_L} L(a_L, \text{target}) \odot \sigma'(z_L)$$

$$W_L = W_L - \delta_L a_{\text{previous}}^T$$

Backward pass previous layers

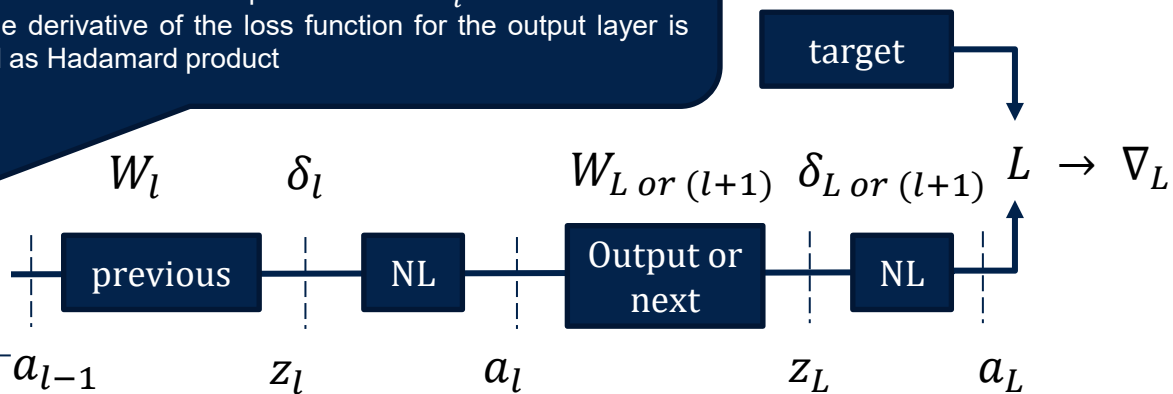
for $l = \text{previous}, \dots$ do

$$\delta_l = W_{l+1}^T \delta_{l+1} \odot \sigma'(z_l)$$

Weight update

$$W_l = W_l - \delta_l a_{l-1}^T$$

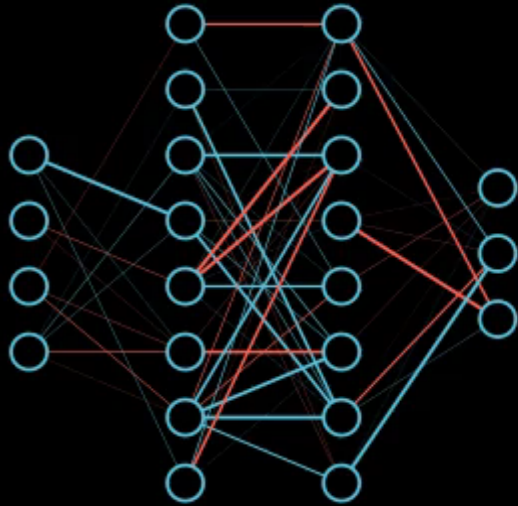
end for



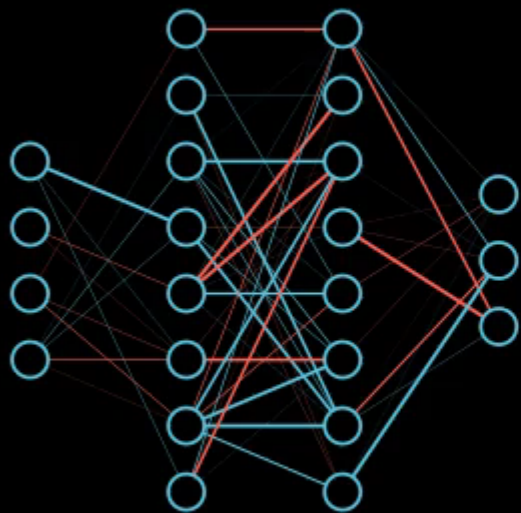
in reverse pipeline order. the derivative of the loss the derivative of previous step is multiplied to the input activations of the previous layer a_{l-1} to compute the variation of the weights that are added to update the weights of the previous layer

GT = ground truth
NL = non-linearity
L = Loss

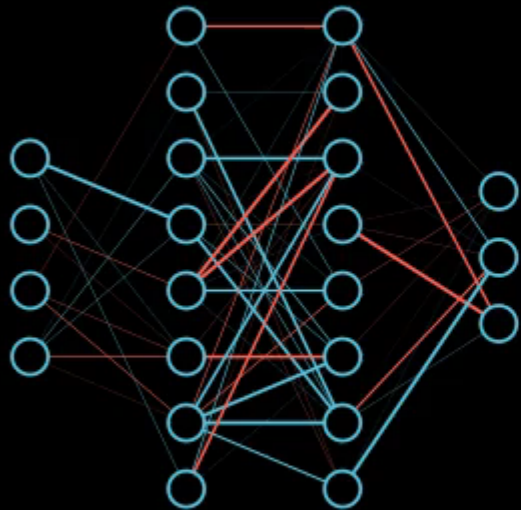
Backpropagation



Backpropagation



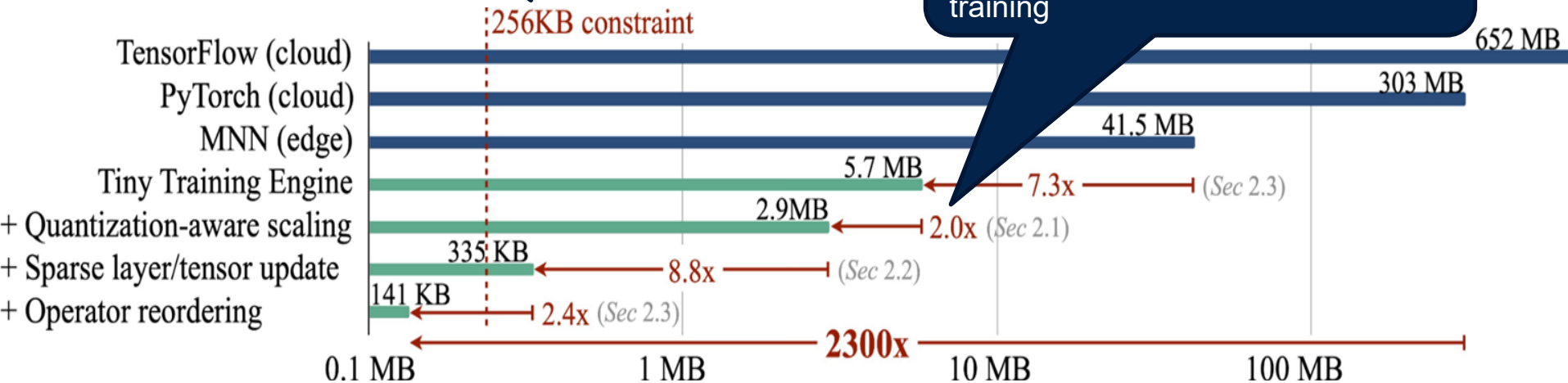
Backpropagation



MCUNetV3²

State of art on ODL
CNN learning in only
256KiB SRAM
Visual Wake Words

QAS (Quantization Aware Scaling):
mitigated the backprop instability due
to the int8 quantization error during
training



Sparse Update:

- updates only **some** parameters of the model.
- these parameters are selected **offline** according to how much they contribute to reduce the error during training.
- it needs to store only **the intermediate activations** of such parameters

+ Sparse layer/tensor update |  335 KB ← 8.8x (Sec 2.2)

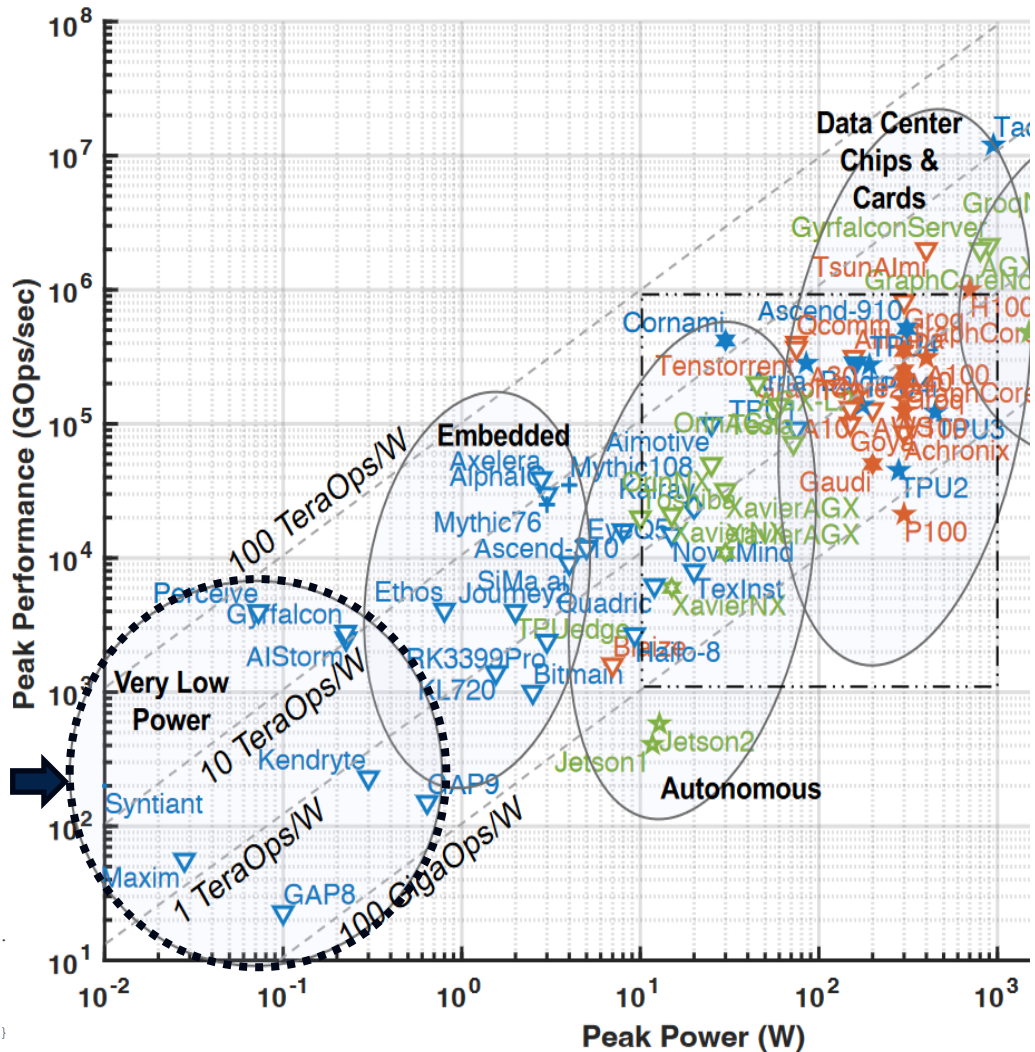
Hardware acceleration

Hardware acceleration is the use of computer hardware designed to perform specific functions more efficiently when compared to software running on a general-purpose central processing unit (CPU). *

Application	Hardware accelerator	Acronym
Computer graphics	<ul style="list-style-type: none"> •General-purpose computing on GPU •CUDA architecture •Ray-tracing hardware 	<ul style="list-style-type: none"> •GPGPU •CUDA •RTX
Digital signal processing	Digital signal processor	DSP
Analog signal processing	<ul style="list-style-type: none"> •Field-programmable analog arrayField-programmable RF 	•FPAAFPRF
Sound processing	Sound card and sound card mixer	N/A
Computer networking on a chip	<ul style="list-style-type: none"> •Network processor and network interface controller Network on a chip 	•NPU and NICNoC
Cryptography Encryption Attack Random number generation	<ul style="list-style-type: none"> •Cryptographic accelerator and secure cryptoprocessorHardware-based encryption •Custom hardware attack •Hardware random number generator 	N/A
Artificial intelligence Machine vision/computer vision Neural networks Brain simulation	<ul style="list-style-type: none"> •AI acceleratorVision processing unit •Physical neural network •Neuromorphic engineering 	<ul style="list-style-type: none"> •N/AVPU •PNN •N/A
Multilinear algebra	Tensor processing unit	TPU
Physics simulation	Physics processing unit	PPU
Regular expressions ^[16]	Regular expression coprocessor	N/A
Data compression ^[17]	Data compression accelerator	N/A
In-memory processing	Network on a chip and Systolic array	NoC; N/A
Data processing	Data processing unit	DPU
Any computing task	<ul style="list-style-type: none"> •Computer hardwareField-programmable gate arrays^[18] •Application-specific integrated circuits^[18] •Complex programmable logic devices •Systems-on-Chip 	<ul style="list-style-type: none"> •HW (sometimes)FPGA •ASIC •CPLD •SoC

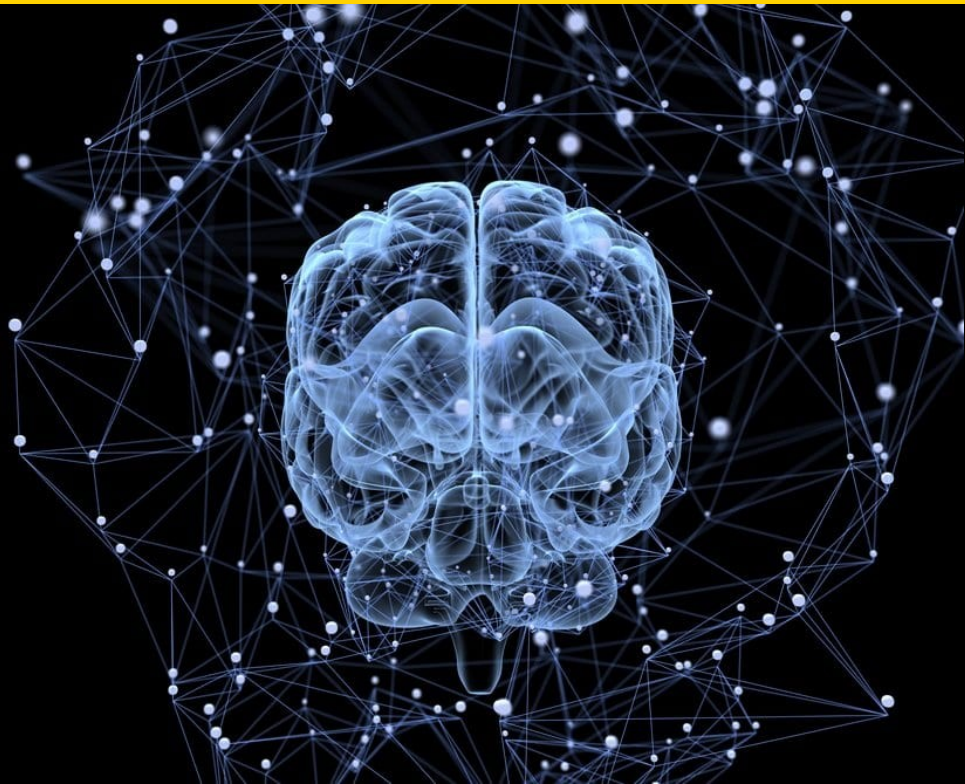
Edge AI drives **1** hardware for on-device machine learning inference.

Tiny ML hardware (today) is optimized for forward workloads **2**



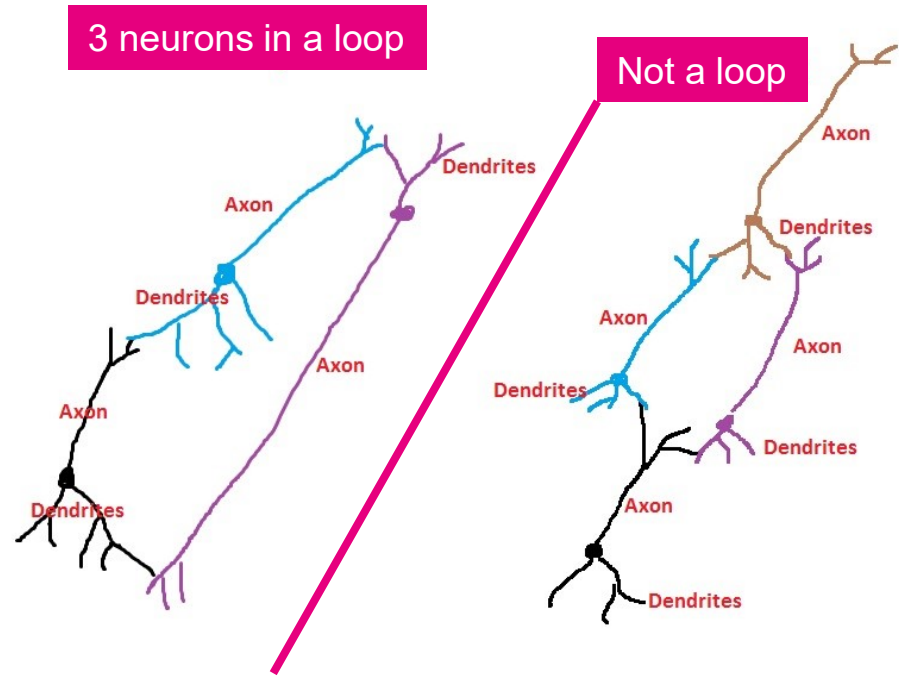
@article{Reuther2019SurveyAB, title={Survey and Benchmarking of Machine Learning Accelerators}, author={A. Reuther and Peter Michaleas and Michael Jones and Vijay Gadepally and Siddharth Samsi and Jeremy Kepner}, journal={2019 IEEE High Performance Extreme Computing Conference (HPEC)}, year={2019}, pages={1-9}, url={https://api.semanticscholar.org/CorpusID:201668230} }

Forward-Forward and PEPITA



Learning: biological plausibility

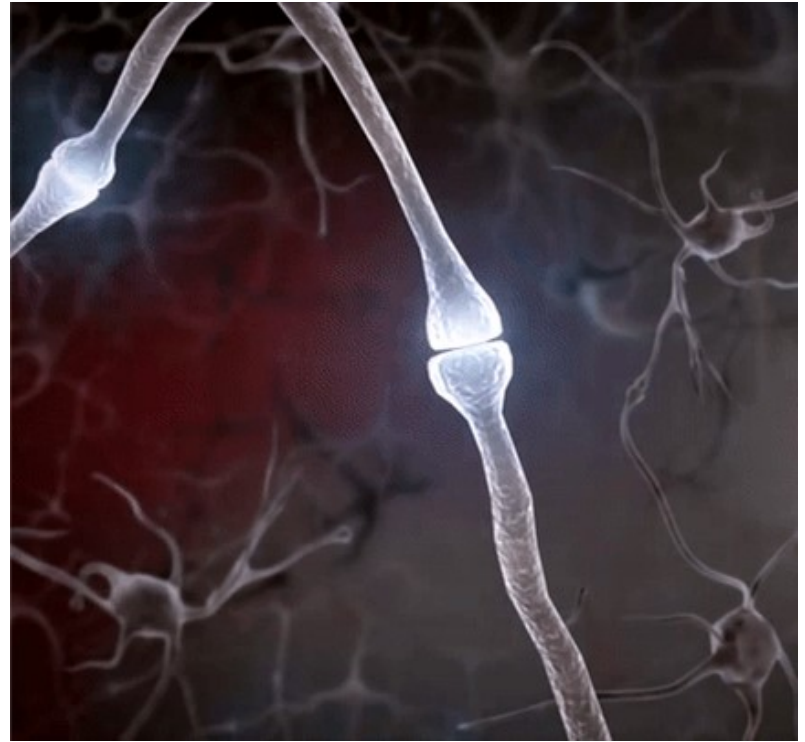
- ❑ The brain **learns** by modifying the synaptic individual connections between neurons³
- ❑ It's **not known how** the single modifications are coordinated to achieve a global's goal
- ❑ **Loop-based** neuron circuits seems used to get error signals and credits (i.e. how much each synapse contributes to the error) assigned to other synapsis of neurons



³T. Lillicrap, A. Santoro, L. Marris, C. Akerman, and G. Hinton, "Backpropagation and the brain," *Nature Reviews Neuroscience*, vol. 21, no. 6, p. 335–346, 2020.

Backpropagation vs bio-plausibility?

Our brain does not use **backpropagation**⁴



⁴Crick, "The recent excitement about neural networks," *Nature*, vol. 337, no. 6203, p. 129—132, January 1989. [Online]. Available: <https://doi.org/10.1038/337129a0>

Backpropagation vs bio-plausibility?

Our brain does not use **backpropagation**⁴

1. No weight symmetry

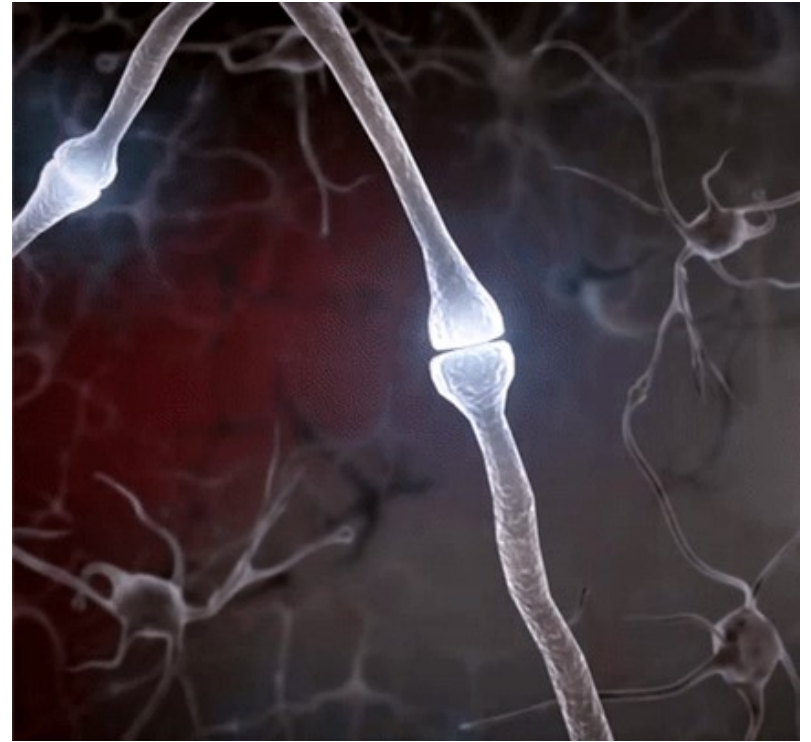
- a) error is not projected back using the same weights of the forward pass



Backpropagation vs bio-plausibility?

Our brain does not use **backpropagation**⁴

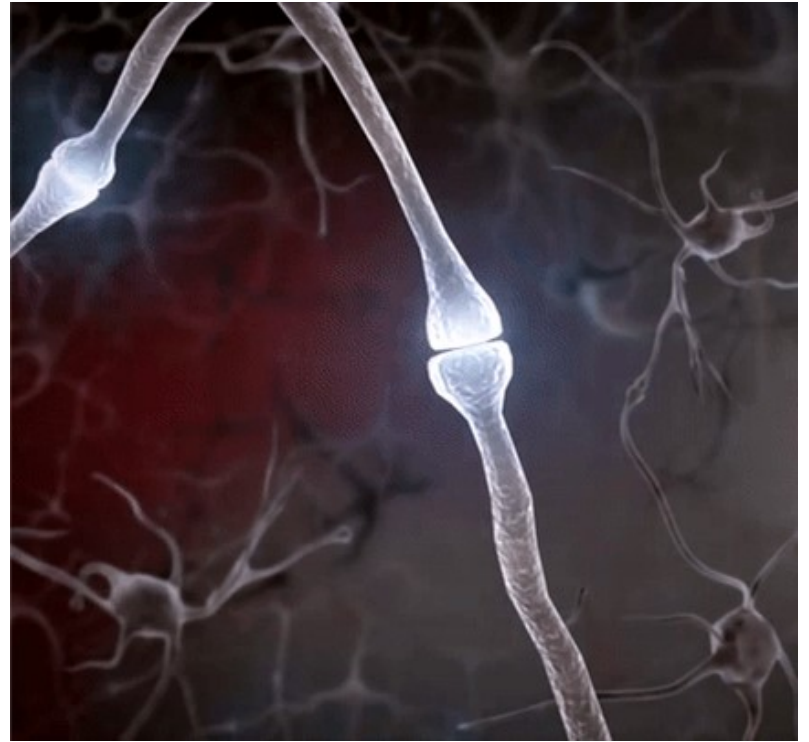
1. No weight symmetry
2. No neural activity freeze
 - a) intermediate activations are not stored



Backpropagation vs bio-plausibility?

Our brain does not use **backpropagation**⁴

1. No weight symmetry
2. No neural activity freeze
3. No locality of the loss function
 - a) Neurons do exchange error signals and credits within loops
 - b) Synapsis learn from local signals










Backpropagation vs bio-plausibility?

Our brain does not use **backpropagation**⁴










1. No weight symmetry
2. No neural activity freeze
3. No locality of the loss function
4. No Update-locking
 - a) No need to wait the end of the backward pass to update the weights of the layers
 - b) There is no backward pass



A lot of work in 30 years...

	Biologically plausible method	Citation
	DRTP DRTP Direct Random Target Projection	C. Frenkel, M. Lefebvre and D. Bol, "Learning without feedback: Fixed Random Learning Signals Allow for Feedforward Training of Deep Neural Networks," <i>Frontiers in Neuroscience</i> , vol. 15, no. 629892, 2021. doi: 10.3389/fnins.2021.629892
	GEVB GEVB Global Error Vector Broadcasting	Clark, David, L. F. Abbott, and SueYeon Chung. "Credit assignment through broadcasting a global error vector." <i>Advances in Neural Information Processing Systems</i> 34 (2021): 10053-10066.
	WM WM Weight Mirroring	Akrout, Mohamed, et al. "Deep learning without weight transport." <i>Advances in neural information processing systems</i> 32 (2019).
	FA FA Feedback Alignment	Lillicrap, Timothy P., et al. "Random feedback weights support learning in deep neural networks." <i>arXiv preprint arXiv:1411.0247</i> (2014).
	DFA DFA Direct Feedback Alignment	Nøkland, Arild. "Direct feedback alignment provides learning in deep neural networks." <i>Advances in neural information processing systems</i> 29 (2016).
	DFC DFC Deep Feedback Control	Meulemans, Alexander, et al. "Credit assignment in neural networks through deep feedback control." <i>Advances in Neural Information Processing Systems</i> 34 (2021): 4674-4687.
	CLAPP CLAPP Contrastive Local And Predictive Plasticity	Illing, Bernd, et al. "Local plasticity rules can learn deep representations using self-supervised contrastive predictions." <i>Advances in Neural Information Processing Systems</i> 34 (2021): 30365-30379.

A lot of work in 30 years...

	Biologically plausible method	Citation
 DRTP	DRTP Direct Random Target Projection	C. Frenkel, M. Lefebvre and D. Bol, "Learning without feedback: Fixed Random Learning Signals Allow for Feedforward Training of Deep Neural Networks," <i>Frontiers in Neuroscience</i> , vol. 15, no. 629892, 2021. doi: 10.3389/fnins.2021.629892
 GEVB	GEVB Global Error Vector Broadcasting	Clark, David, L. F. Abbott, and SueYeon Chung. "Credit assignment through broadcasting a global error vector." <i>Advances in Neural Information Processing Systems</i> 34 (2021): 10053-10066.
 WM	WM Weight Mirroring	Akrout, Mohamed, et al. "Deep learning without weight transport." <i>Advances in neural information processing systems</i> 32 (2019).
 FA	FA Feedback Alignment	Lillicrap, Timothy P., et al. "Random feedback weights support learning in deep neural networks." <i>arXiv preprint arXiv:1411.0247</i> (2014).
 DFA	DFA Direct Feedback Alignment	Nøkland, Arild. "Direct feedback alignment provides learning in deep neural networks." <i>Advances in neural information processing systems</i> 29 (2016).
 DFC	DFC Deep Feedback Control	Meulemans, Alexander, et al. "Credit assignment in neural networks through deep feedback control." <i>Advances in Neural Information Processing Systems</i> 34 (2021): 4674-4687.
 CLAPP	CLAPP Contrastive Local And Predictive Plasticity	Illing, Bernd, et al. "Local plasticity rules can learn deep representations using self-supervised contrastive predictions." <i>Advances in Neural Information Processing Systems</i> 34 (2021): 30365-30379.
 FF	Forward-Forward	Hinton, Geoffrey. "The forward-forward algorithm: Some preliminary investigations." <i>arXiv preprint arXiv:2212.13345</i> (2022).
 PEPITA	Present the Error to Perturb the Input to modulate Activity	Dellaferrera, Giorgia, and Gabriel Kreiman. "Error-driven input modulation: solving the credit assignment problem without a backward pass." <i>International Conference on Machine Learning</i> . PMLR, 2022.

Positive data

Negative data

Forward-Forward⁵

Algorithm 1 Forward-Forward

Given: Positive input (x_{pos}) and Negative input (x_{neg})

Positive Pass

```

 $a_0 = x_{pos}$ 
for  $\ell = 1, \dots, L$  do
   $a_\ell = ReLU(W_\ell a_{\ell-1})$ 
  Weight update
   $W_\ell = W_\ell + 2 \frac{\partial \log(p)}{\partial G(a_\ell)} a_\ell \cdot a_{\ell-1}^T$ 
   $a_\ell = norm(\ell)$ 
end for

```

Positive weight update of each layer independently from the next ones (locality)

Forward inference

Activations of each layer are normalized

Negative Pass

```

 $a_0 = x_{neg}$ 
for  $\ell = 1, \dots, L$  do
   $a_\ell = ReLU(W_\ell a_{\ell-1})$ 
  Weight update
   $W_\ell = W_\ell - 2 \frac{\partial \log(p)}{\partial G(a_\ell)} a_\ell \cdot a_{\ell-1}^T$ 
   $a_\ell = norm(a_\ell)$ 
end for

```

Negative weight update of each layer independently from the next ones (locality)

Forward inference

Activations of each layer are normalized

⁵G. Hinton, "The forward-forward algorithm: Some preliminary investigations," ArXiv, vol. abs/2212.13345, 2022.

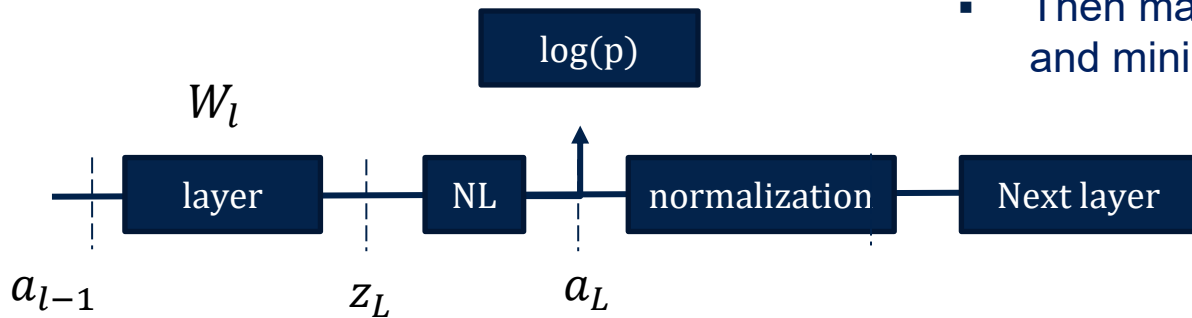
Forward-Forward⁵

- 2 variants: supervised and unsupervised
→ **same learning procedure, different forward procedures**
- **Supervised variant: n passes, for n classes**

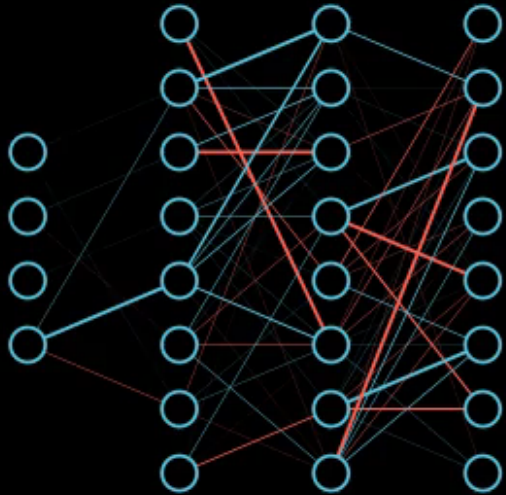
- The **goodness function** $G(a_l)$ is defined at each layer; e.g. it can be the sum of the squared a_l activations
- The probability of being classified as a positive data is

$$p(\text{positive}) = \sigma(G(a_\ell) - \theta)$$

- Then maximize the $\log(p)$ for positive data and minimize it for negative



Forward-Forward



Algorithm 3 PEPITA**Given:** Features(x) and label($target$)**Standard Pass**

$$a_0 = x$$

for $\ell = 1, \dots, L$ **do**

$$a_\ell = \sigma_\ell(W_\ell a_{\ell-1})$$

end for

$$e = a_L - target$$

Modulated pass

$$a_0^{err} = x + Fe$$

for $\ell = 1, \dots, L$ **do**

$$a_\ell^{err} = \sigma_\ell(W_\ell a_{\ell-1}^{err})$$

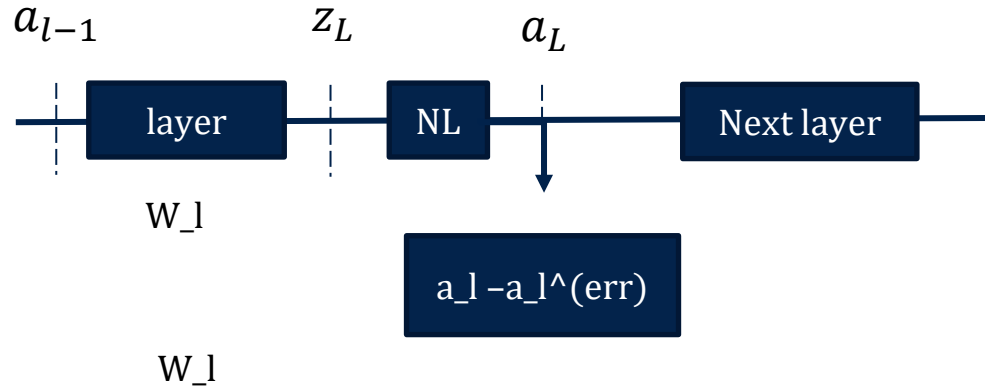
Weight update

$$W_\ell = W_\ell - (a_\ell - a_\ell^{err}) \cdot (a_{\ell-1}^{err})^T$$

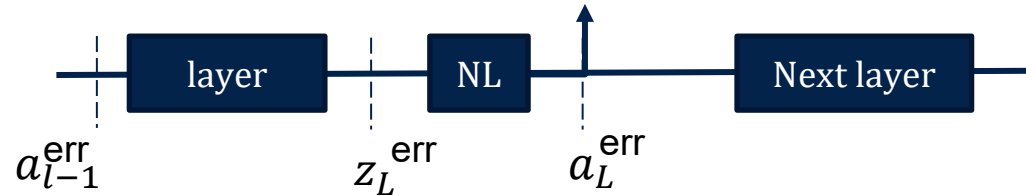
end forForward inference
Store activations #Compute the
error at the outputRandom **matrix F**, with zero
mean and small variance, to
project the **error on the inputs**Forward inference
Store activations #Weights update by using the
forward pass computed activations
(**stored in memory #**) with the
modulated ones

$$W_\ell = W_\ell - (a_\ell - a_\ell^{err}) \cdot (a_{\ell-1}^{err})^T$$

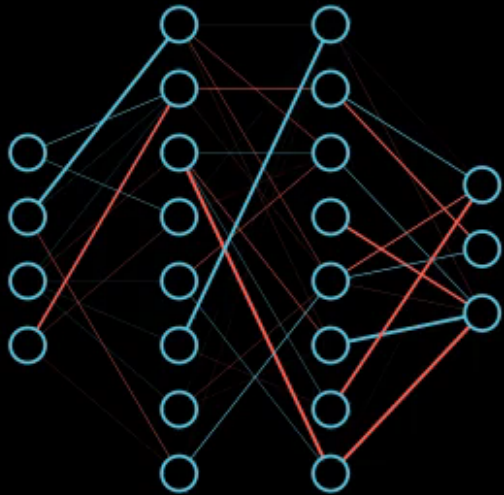
Standard pass



Modulated pass



PEPITA



Research question

At which computational and memory cost would FF and PEPITA learning algorithms compare to BP if applied to MLCommons/Tiny benchmarks ?

Introducing **MEMPEPITA** to not store intermediate activations. Memory savings expected !

Computational complexity and memory footprint of **FF**, **PEPITA** and **MEMPEPITA** for the **MLCommons/Tiny benchmarks**

Algorithm 2 MEMPEPITA**Given:** Features(x) and label($target$)**Standard Pass**

$$a_0 = x$$

for $\ell = 1, \dots, L$ **do**

$$a_\ell = \sigma_\ell(W_\ell a_{\ell-1})$$

end for

$$e = a_L - target$$

Error projection

$$a_0^{err} = x + Fe$$

for $\ell = 1, \dots, L$ **do****Standard pass**

$$a_\ell = \sigma_\ell(W_\ell a_{\ell-1})$$

Modulated pass

$$a_\ell^{err} = \sigma_\ell(W_\ell a_{\ell-1}^{err})$$

Weight update

$$W_\ell = W_\ell - (a_\ell - a_\ell^{err}) \cdot (a_{\ell-1}^{err})^T$$

end for

Forward inference

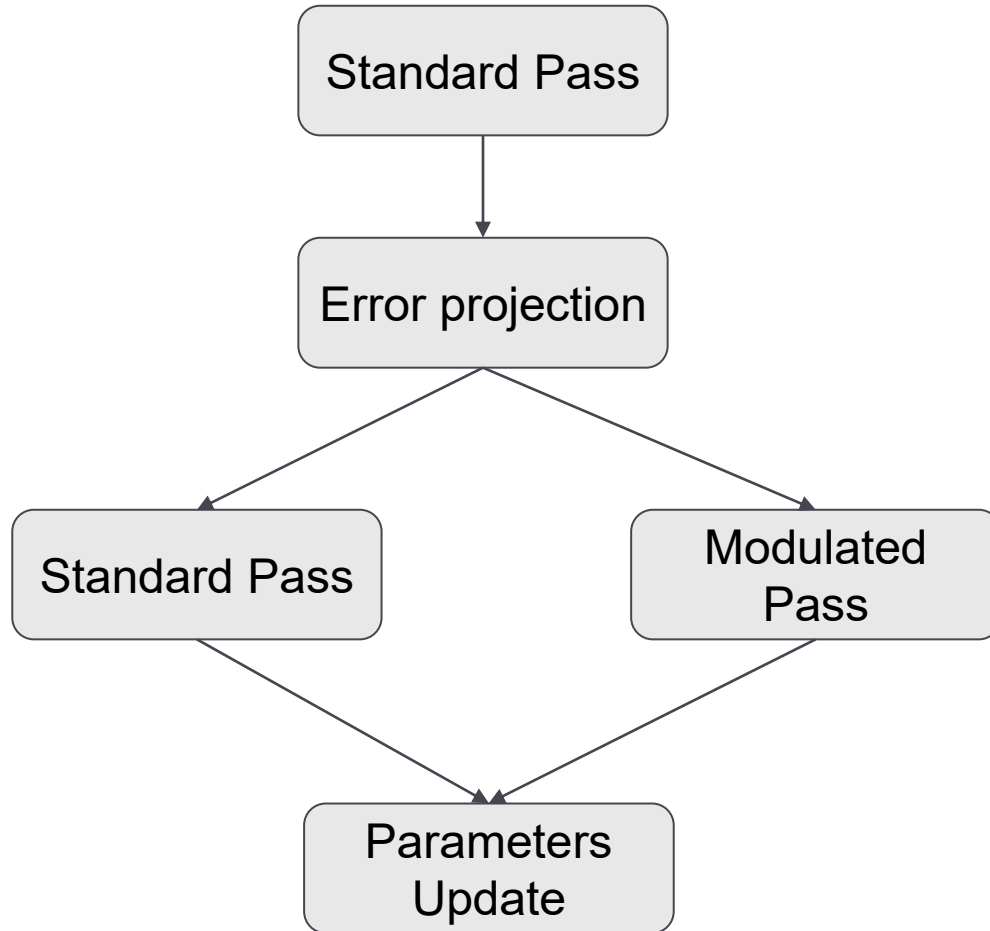
Compute the error at the output

Random **matrix F**, with zero mean and small variance, to project the **error** onto the inputs

Forward inferences thus recomputing activations instead of storing them into memory

Introduces a **second standard pass** which runs simultaneously along with **modulated pass**

Weights update by using the forward pass re-computed activations with the modulated ones



Summary of the learning procedures

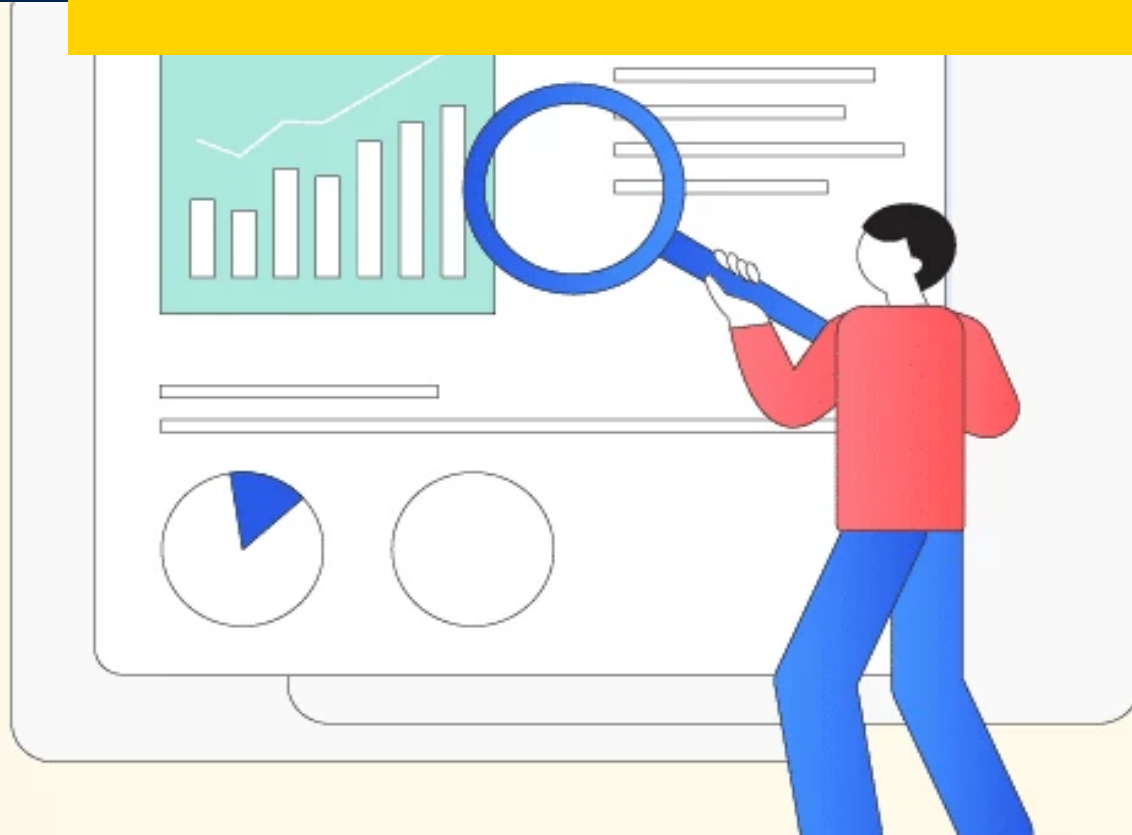
Method	BP (number)	FF (number)	PEP (number)	MPE (number)
Forward passes	1	2	2	3
Backward passes	1	0	0	0
Weight update	1	2	1	1
Loss function	Global	Local	Global	Global
Activations	all	current	all	current

PEP = PEPITA
MPE = MEMPEPITA

Local = loss function per layer
Global = loss function at the output layer

All = all layers
Current = current layer

Methodology



MLCommons/Tiny

Use Case	Description	Dataset	Model	Quality Target (Closed Division)
Audio Wake Words	Small vocabulary keyword spotting	Speech Commands	DS-CNN	90% (Top1)
Visual Wake Words	Image classification (2 classes)	Person Detection	MobileNet	80% (Top1)
Image Classification	32x32 tiny Images Classification (10 classes)	Cifar10	ResNet	85% (Top1)
Anomaly Detection	Detecting anomalies in machine operating sounds	ToyADMOS	Deep AutoEncoder	0.85 (AUC)

ML

- **Commons** framework specifies number of samples and epochs
- Weights, biases, activations represented in INT8
- Softmax layer represented in FLOAT32
- MACCs represented in INT8
- No layer memory overwrite
- Batch normalization not considered.
- Cycles/MACC* and processor's frequency
- Results validated with **STM32Cube.AI Developer Cloud**

MCU deployability

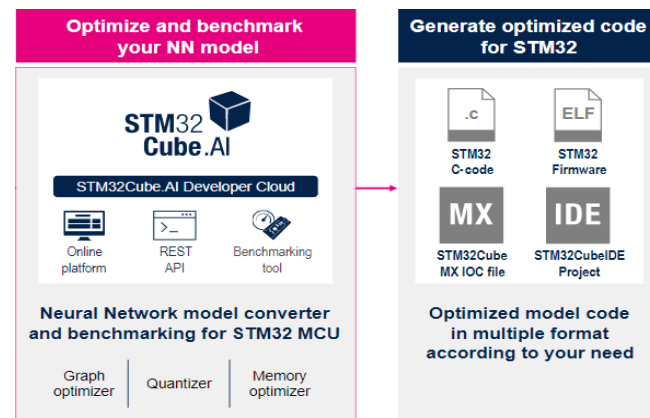
NUCLEO-G474RE
170MHz



STM32H735G-DK
550MHz



<https://stm32ai-cs.st.com/home>



Complexity analysis: Assumptions

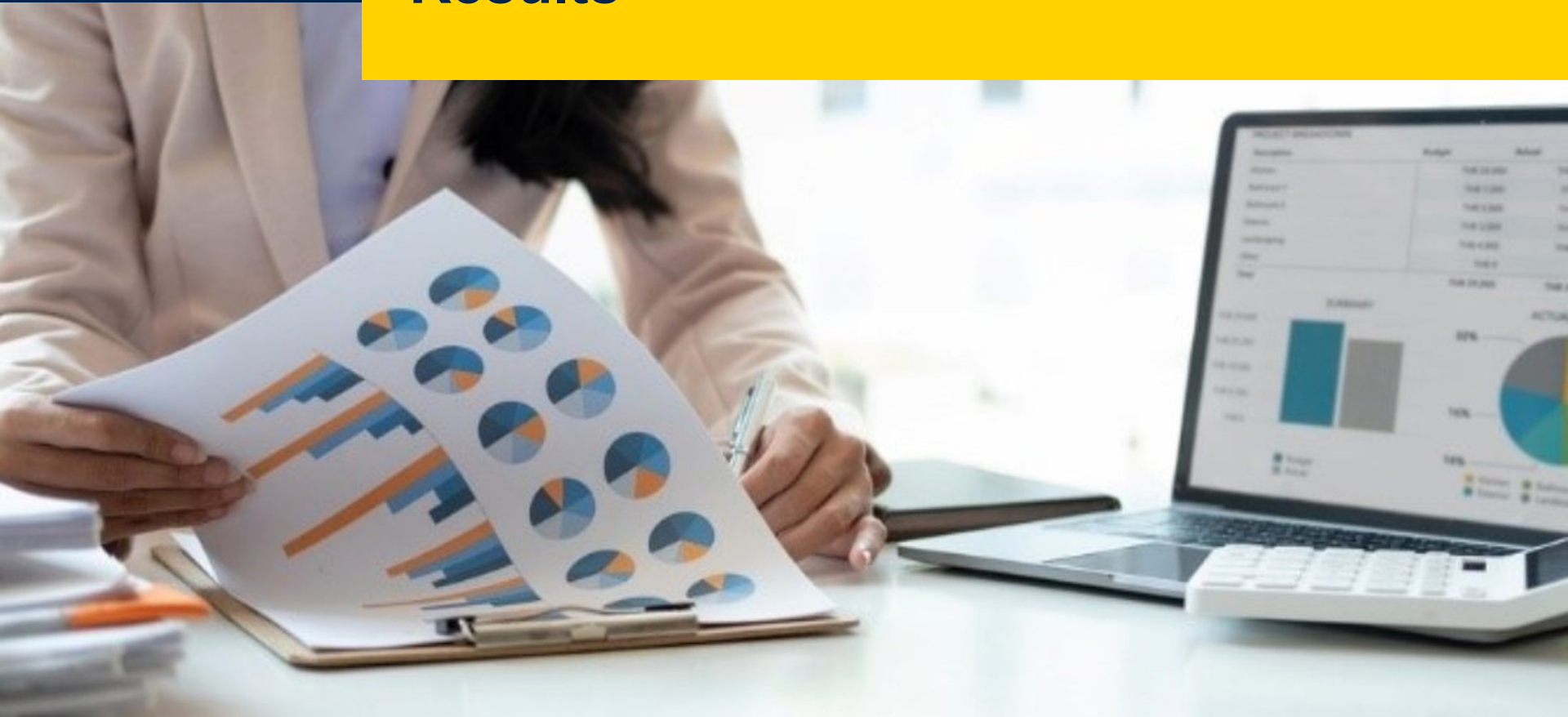
Computations required	Forward pass	Error at output layer	Backward pass	Weight update	Goodness function	Normalization	Error projection
BP	1	1	1	1			
FF	2			2	2	2	
PEP	2	1		1			1
MPE	3	1		1			1

RAM estimation

Learning procedure	Activations RAM during training
BP and PEP	Sum of the activation buffers of all layers
FF	Max value of the sum of the activation buffers of two consecutive layers + the input sample
MPE	Max value of the sum of the activation buffers of two consecutive layers + the largest activation buffer between these two layers

Total RAM = Activations RAM + footprint (weights + biases)

Results



Learning procedure: the analysis

Model/Dataset		DS-CNN/SC/KWS				MobileNet/VWW				ResNet/Cifar10				AE/ToyADMOS/Anomaly Detection			
Learning method		BP	FF	PEP	MPE	BP	FF	PEP	MPE	BP	FF	PEP	MPE	BP	FF	PEP	MPE
Training	MAC C (M)	7.7	+43%	+4%	+39%	22	+41%	+3%	+37%	37	+36%	+1%	+35%	0.7	+50%	+69%	+106%
	ACT (KIB)	+253%	21	+253%	+37%	+213%	83	+213%	+11%	+118%	52	+118%	+57%	+64%	1.4	+64%	1.4
	PAR/ACT	0.3	1.1	0.3	0.8	0.8	2.5	0.8	2.3	0.7	1.5	0.7	0.9	115	189	115	189
	RAM (KIB)	+120%	43	+120%	+17%	+60%	294	+60%	+3%	+47%	130	+47%	+23%	+0.3%	267	+0.3%	267

averaged

Architecture		CNN			FC(AE)		
Learning method		FF	PEP	MPE	FF	PEP	MPE
Train	MACC	+40%	+3%	+37%	+50%	+69%	+106%
	ACT	-65%	0%	-51%	-39%	0%	-39%
	RAM	-41%	0%	-33%	-0.3%	0%	-0.3%

Inference procedure: the analysis

Model/Dataset		DS-CNN/SC				MobileNet/VWW				ResNet/Cifar10				AE/ToyADMOS			
Inference method		BP	FF	PEP	MPE	BP	FF	PEP	MPE	BP	FF	PEP	MPE	BP	FF	PEP	MPE
Inference	MACC (M)	3	+1167%	3	3	8	+113%	8	8	13	+914%	13	13	0.3	0.3	0.3	0.3
	RAM (KiB)	20	+2.4%	20	20	55	+50%	55	55	49	+6%	49	49	0.8	0.8	0.8	0.8
	ROM (KiB)	22.604				210.85				77.706				265.864			

averaged

		CNN			FC(AE)		
Inference method		FF	PEP	MPE	FF	PEP	MPE
Inference	MACC	+731%	0%	0%	+1%	0%	0%
	RAM	+20%	0%	0%	0%	0%	0%

Latency per input sample on MCUs

H7 = STM32H735G-DK @ 550 MHz
 G4 = NUCLEO-G474RE @ 170 MHz

Model/Dataset		DS-CNN/SC				MobileNet/VWW				ResNet/Cifar10				AE/ToyADMOS			
Learning method		BP	FF	PEP	MPE	BP	FF	PEP	MPE	BP	FF	PEP	MPE	BP	FF	PEP	MPE
Training (ms)	H7	42	+43%	+4%	+39%	122	+41%	+3%	+37%	202	+36%	+1%	+35%	4	+50%	+69%	+106%
	G4	203	+43%	+4%	+39%	592	+41%	+3%	+37%	982	+36%	+1%	+35%	19	+50%	+69%	+106%
Inference (ms)	H7	14.5	+1165%	14.5	14.5	42	+113%	42	42	68	+914%	68	68	1.4	+1%	1.4	1.4
	G4	71	+1165%	71	71	202	+113%	202	202	332	+914%	332	332	7	+1%	7	7

- **LEARNING → FF and MEMPEPITA**
 - reduced activations (ACT) storage on average **40% to 65% (w.r.t. BP)**,
 - increased MACCs **40% to 100%**
 - **PEPITA** (same memory as BP) increased MACCS (CNN) **3%**, (FC) **69%**.
- Total **RAM reduction is noticeable** if the topology has **low parameters/activations** (FC vs CNN). → e.g. for DS-CNN is **0.3** and RAM reduction is around **-100%**
- **INFERENCE → MEMPEPITA, PEPITA and BP** featured **1 forward** pass, while supervised **FF** adds **2-3x** more computation due to **N forward** passes, for N classes

REMEMBER

About training → MEMPEPITA reduced total RAM, (CNN) 33%, (FC) 0.3%, at the expense of a third more MACCs.

Inference complexity was unchanged

<https://github.com/fabrizioaymone/suitability-of-Forward-Forward-and-PEPITA-learning>

main 1 branch 0 tags

Go to file Add file Code

fabrizioaymone Add files via upload 5f05fc7 4 days ago 40 commits

figures	Add files via upload	last week
LICENSE	Create LICENSE	3 weeks ago
README.md	Update README.md	last week
analysis.ods	Add files via upload	4 days ago

README.md

Suitability of Forward-Forward and PEPITA Learning to MLCommons-Tiny benchmarks

This repository contains the spreadsheet("analysis.ods") of the quantitative analysis performed for the paper "Suitability of Forward-Forward and PEPITA Learning to MLCommons-Tiny benchmarks".

Brief description

The objective of the analysis is to evaluate the performances in terms of memory usage and complexity of a learning procedure X on a certain model Y tackling a task T on a dataset D. The learning procedures are Backpropagation (BP), Forward-Forward (FF), PEPITA (PEP), MEMPEPITA (MPE). The models evaluated were respectively named DS-CNN, MobileNet, ResNet and AutoEncoder (AE). The datasets used were Speech Commands (SC), Visual Wake

About

This repository contains the spreadsheet of the quantitative analysis performed for the paper "Suitability of Forward-Forward and PEPITA Learning to MLCommons-Tiny benchmarks".

backpropagation on-device-deep-learning
mlcommons tiny-devices forward-forward
pepita

Readme
View license
0 stars
1 watching
0 forks

Releases

No releases published
[Create a new release](#)

Packages

No packages published
[Publish your first package](#)

Q&A



Further questions ? Please contact: danilo.pau@st.com

Our technology starts with You



Find out more at www.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.

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



life.augmented



Copyright Notice

This multimedia file is copyright © 2023 by tinyML Foundation. All rights reserved. It may not be duplicated or distributed in any form without prior written approval.

tinyML[®] is a registered trademark of the tinyML Foundation.

www.tinyml.org



Copyright Notice

This presentation in this publication was presented as a tinyML® 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