

# tinyML<sup>®</sup> Research Symposium

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

March 27, 2023



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

Symposium  
March 27<sup>th</sup>, 2023  
San Francisco, USA



life.augmented

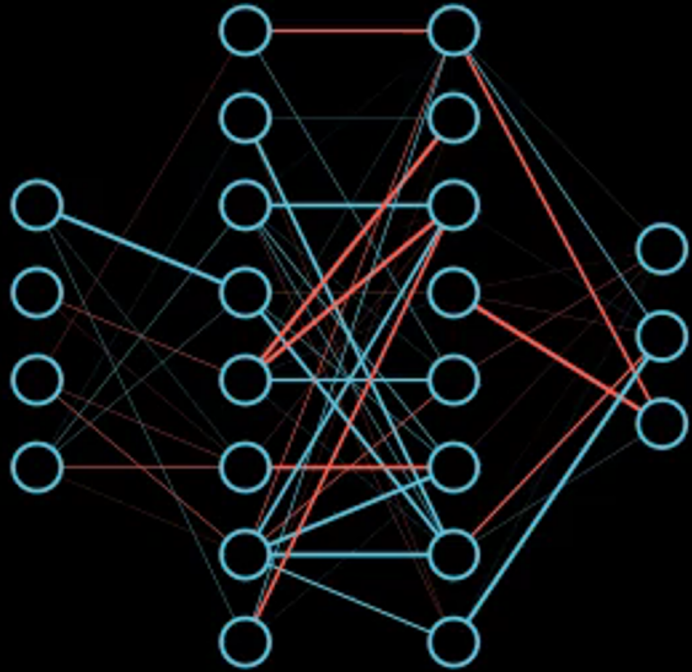
# TinyRCE: Forward Learning Under Tiny Constraints

Danilo Pau, Prem Kumar Ambrose

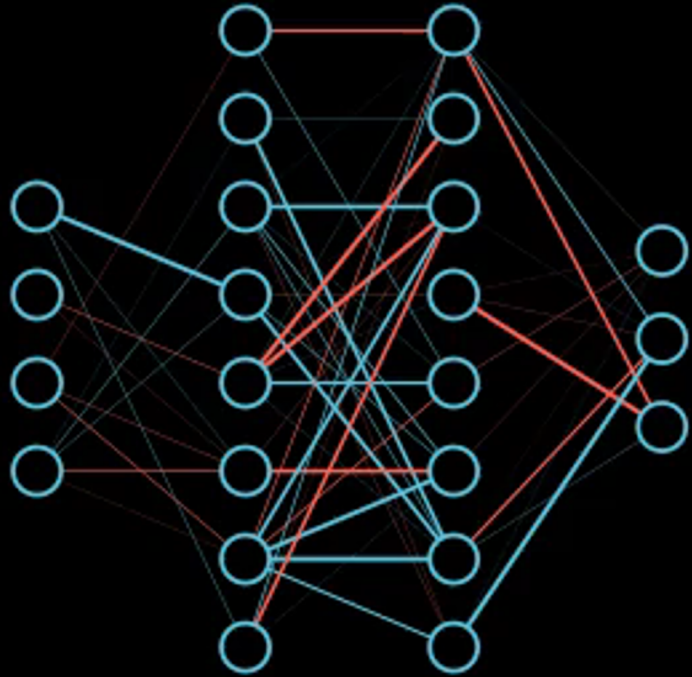
System Research and Applications

STMicroelectronics, Agrate Brianza, Italy

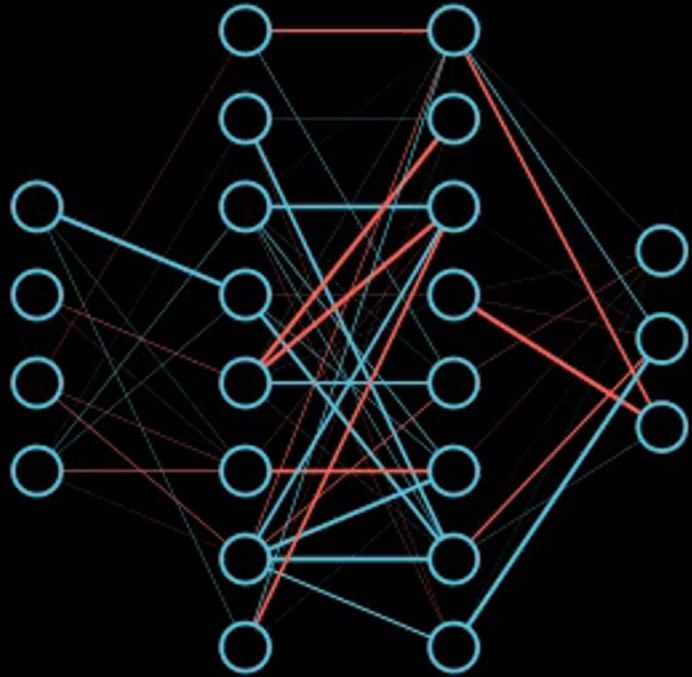
# Backpropagation



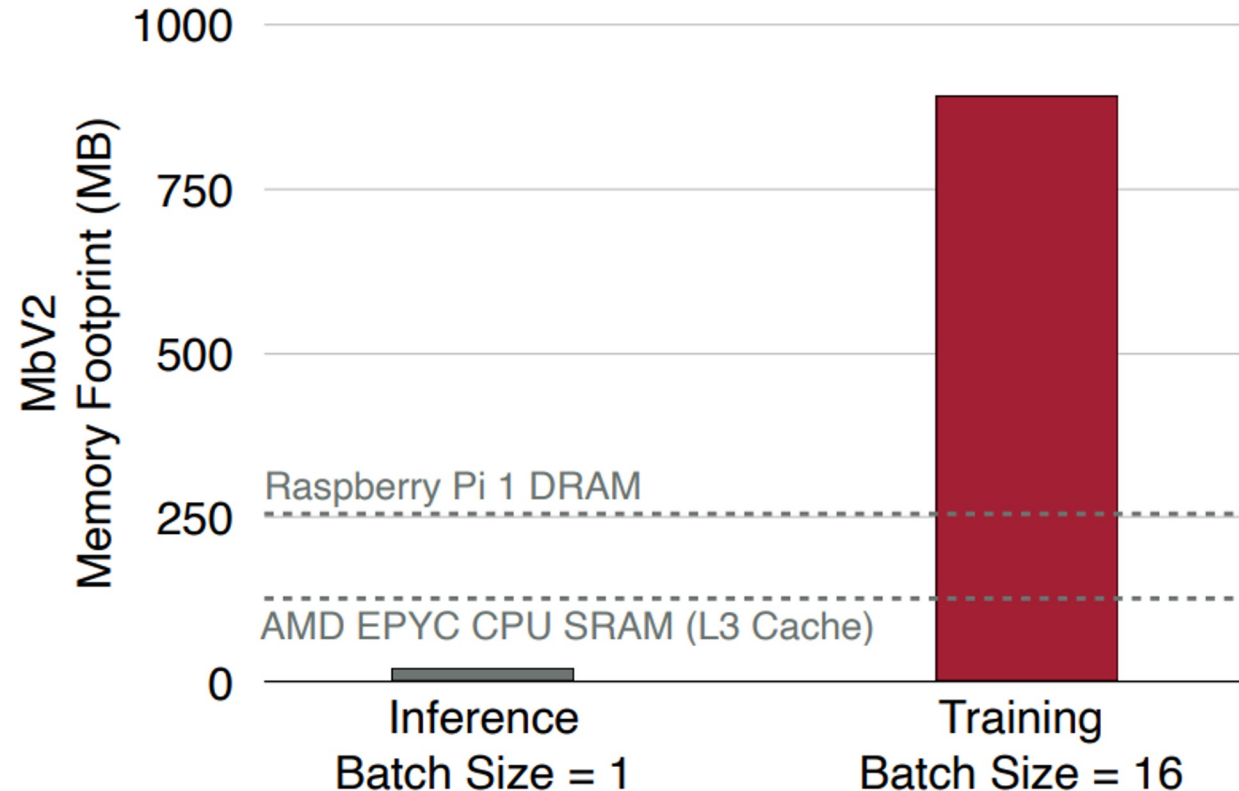
# Backpropagation



# Backpropagation

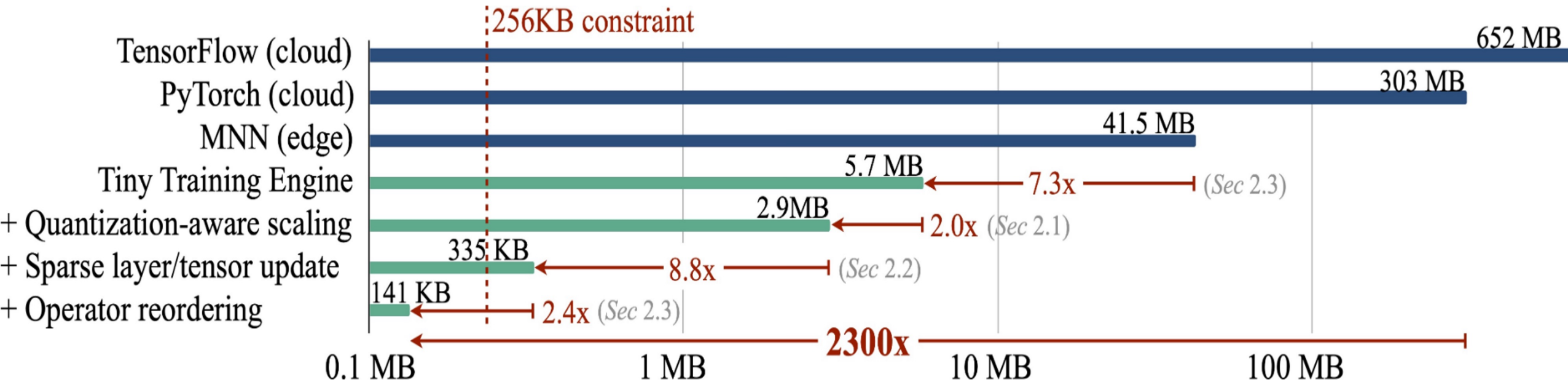


# Reduce activations, not trainable parameters for efficient on-device learning<sup>1</sup>




<sup>1</sup>H. Cai, C. Gan, L. Zhu, and S. Han, "Tinytl: Reduce memory, not parameters for efficient on-device learning," in Advances in Neural Information Processing Systems, vol. 33, 2020

# MCUNetV3<sup>2</sup>



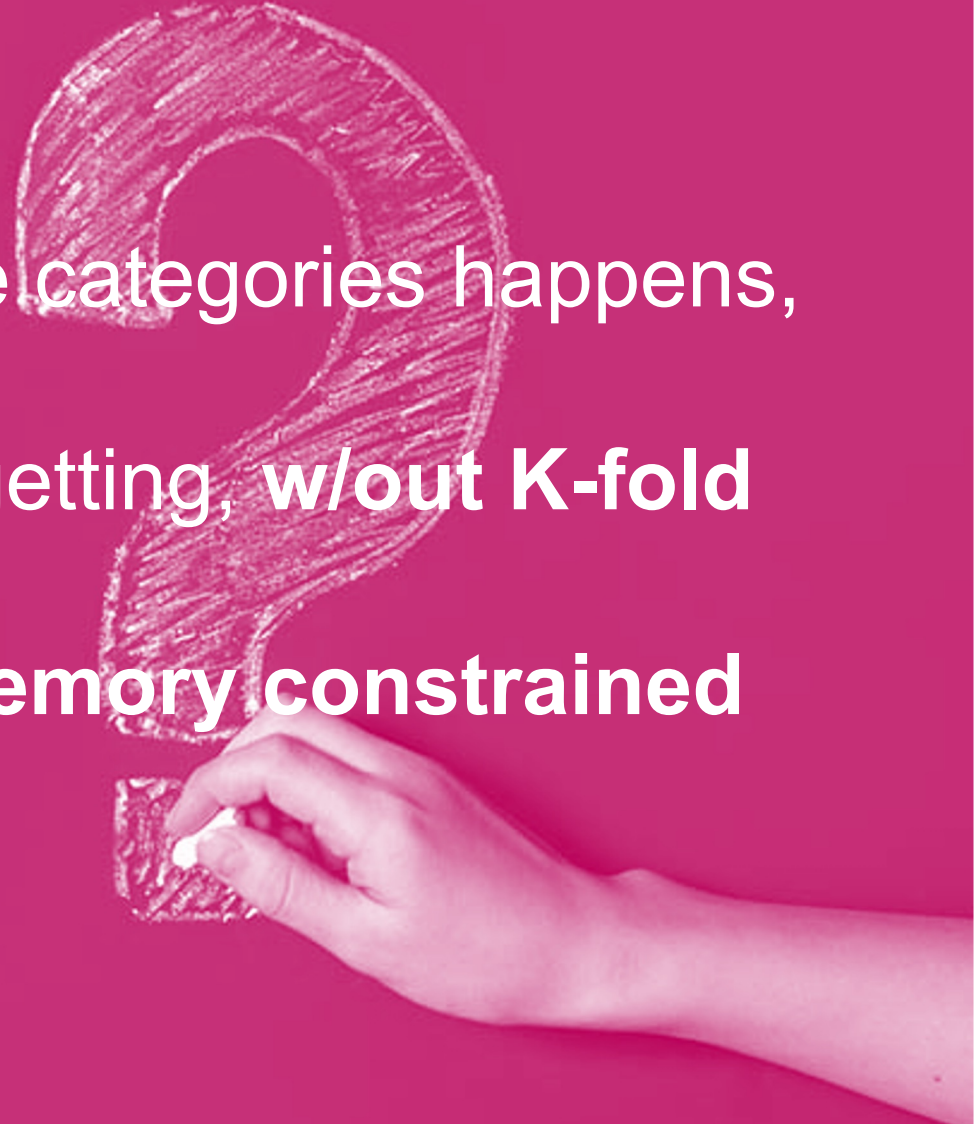
<sup>2</sup>J. Lin, L. Zhu, W.-M. Chen, W.-C. Wang, C. Gan, and S. Han, "On-device training under 256kb memory," 2022. [Online]. Available: <https://arxiv.org/abs/2206.15472>

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

<sup>2</sup>J. Lin, L. Zhu, W.-M. Chen, W.-C. Wang, C. Gan, and S. Han, "On-device training under 256kb memory," 2022. [Online]. Available: <https://arxiv.org/abs/2206.15472>



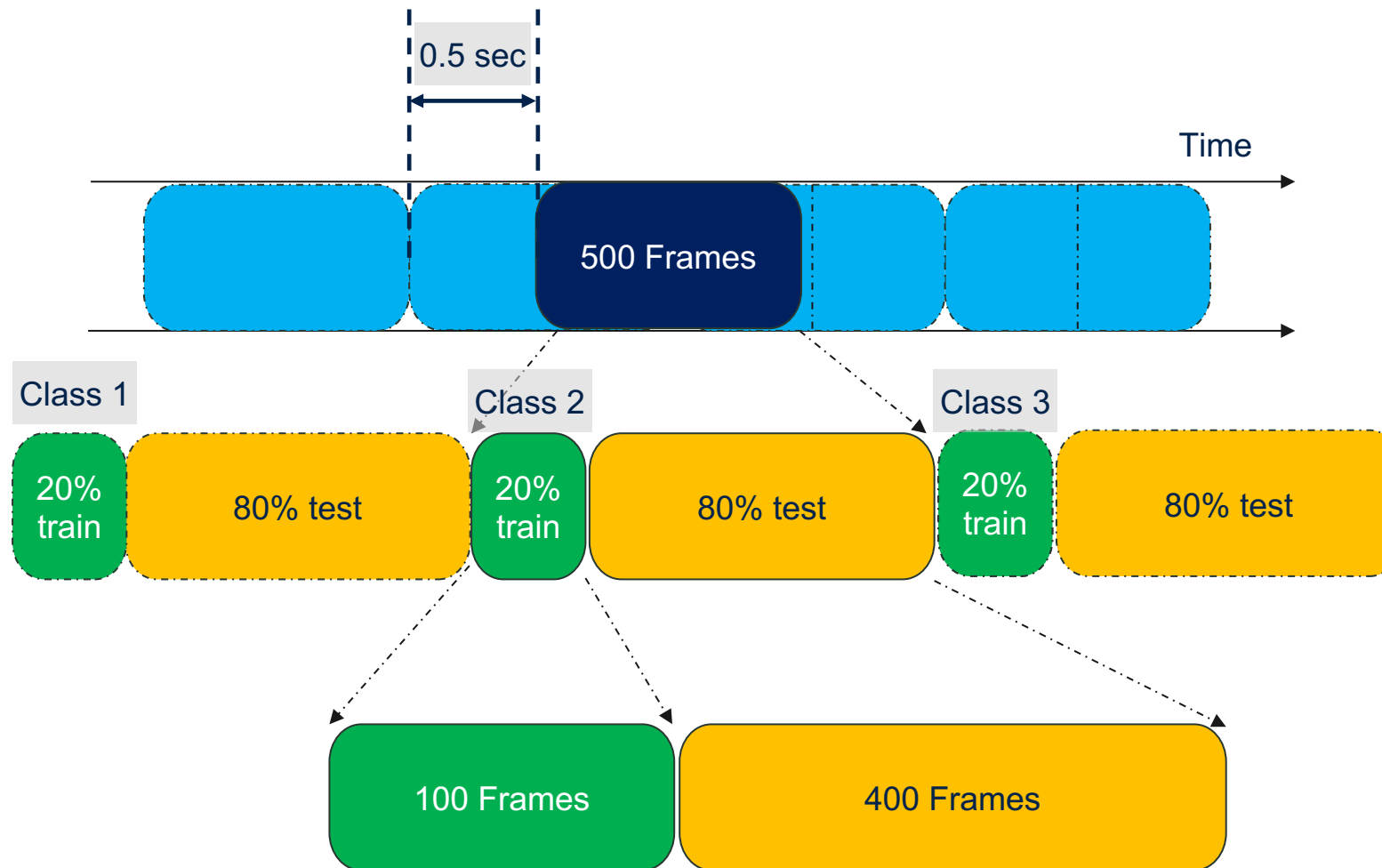
Can **incremental learning** of **multiple** categories happens,  
totally **on-line** w/out catastrophic forgetting, **w/out K-fold**  
**back-prop** and be deployable on **memory constrained**  
devices?

A hand is shown in the bottom right corner, holding a piece of white chalk and drawing a large question mark on a dark green chalkboard. The question mark is drawn with multiple overlapping strokes, giving it a textured, hand-drawn appearance. The hand is positioned as if it has just finished or is in the process of finishing the drawing.

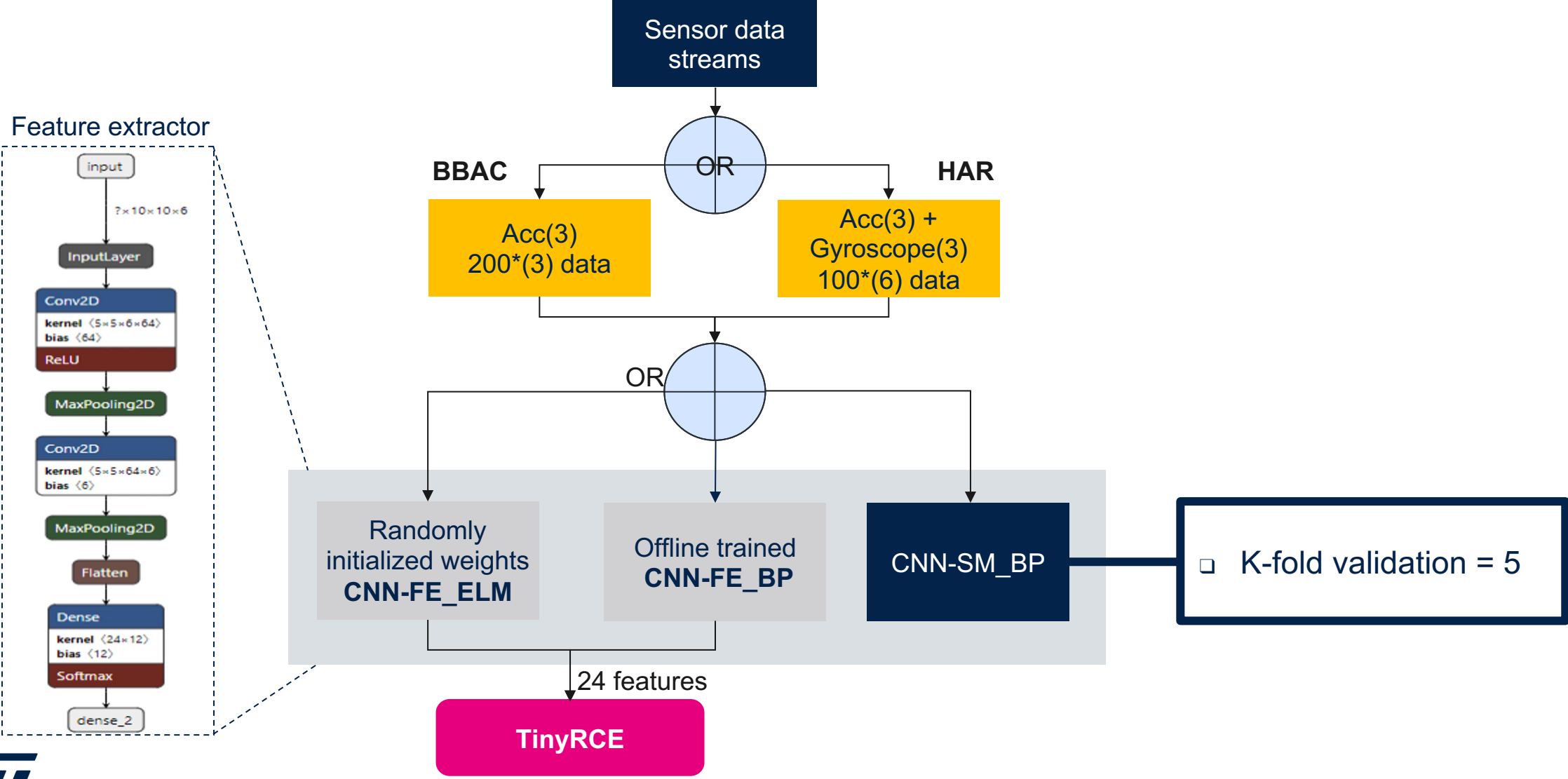
# Case studies

	SHL	PAMAP2	CWRU
Sampling freq [Hz]	100	100	12000
Type	HAR	HAR	BBAC
Sensor data	Accelerometer (3 axis) + gyroscope (3 axis)	Accelerometer (3 axis) + gyroscope (3 axis)	Accelerometer (3 axis)
Number of classes	8	12	10
Number of users	3	8	1
Frame size	1 sec (6 * 100)	1 sec (6 * 100)	16.6ms (3*200)
References	H. Gjoreski et al., "The University of Sussex-Huawei Locomotion and Transportation Dataset for Multimodal Analytics With Mobile Devices," in IEEE Access, vol. 6, pp. 42592-42604, 2018, doi: 10.1109/ACCESS.2018.2858933.	A. Reiss and D. Stricker, "Introducing a New Benchmarked Dataset for Activity Monitoring," 2012 16th International Symposium on Wearable Computers, Newcastle, UK, 2012, pp. 108-109, doi: 10.1109/ISWC.2012.13.	<a href="https://engineering.case.edu/beeringsdatacenter/download-data-file">https://engineering.case.edu/beeringsdatacenter/download-data-file</a>

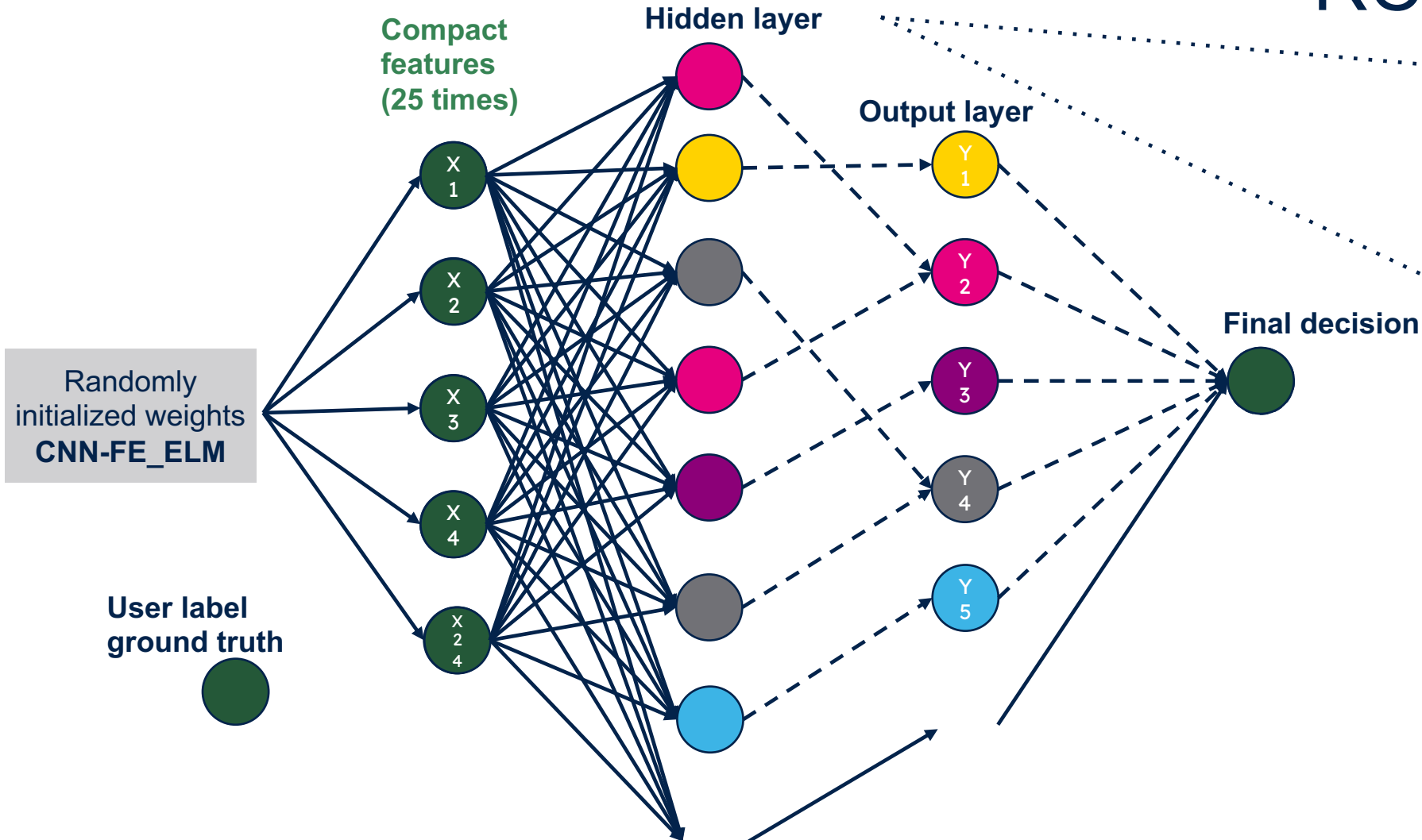
# Interleaved incremental learning and testing



# Proposed framework



# RCE while learning



**Overlapped feature space with input copies stored (like K-NN)**

- Radius keeps changing
- Ambiguous feature assignments
- Fully parallel distance measurements
- No default neural allocation rate control

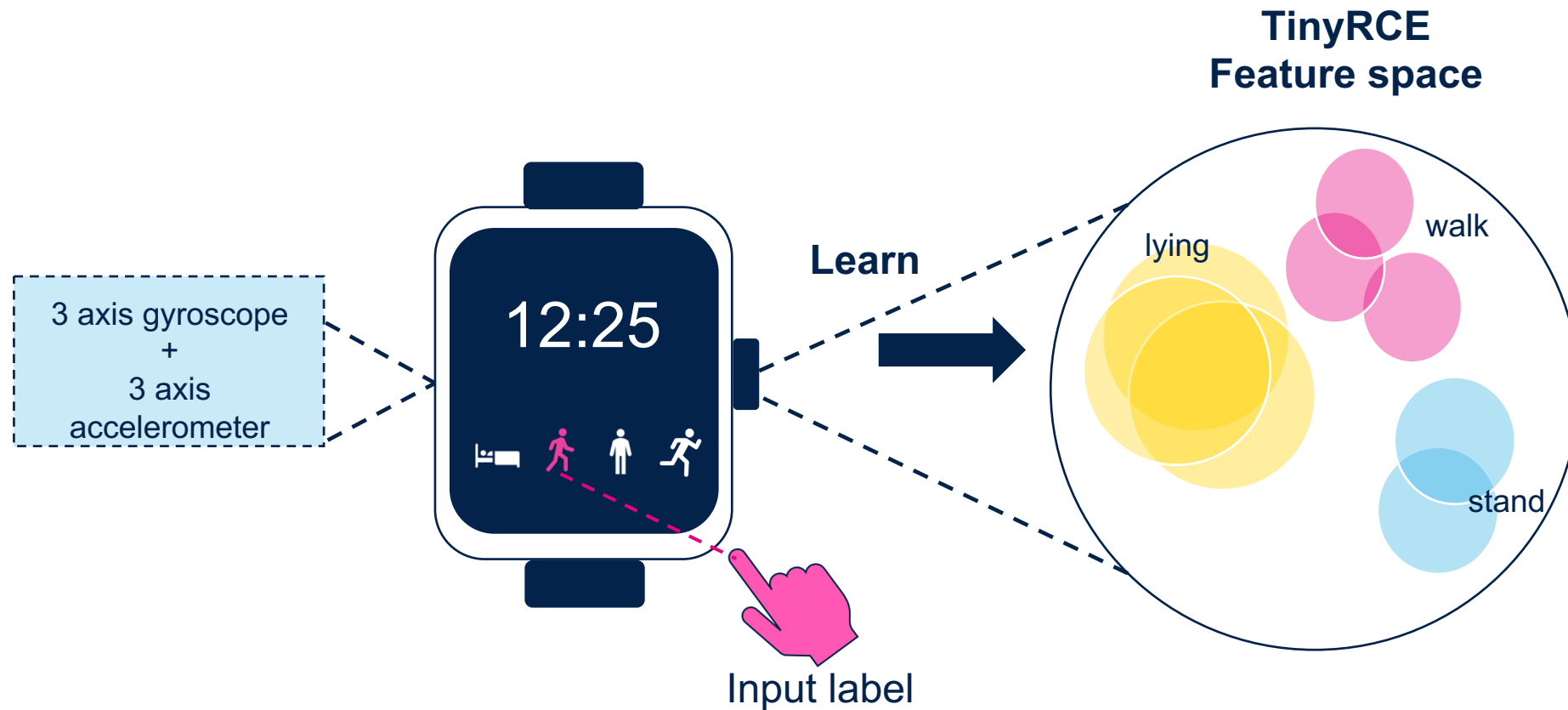
[35] M. J. Hudak, "RCE networks: an experimental investigation," IJCNN-91- Seattle International Joint Conference on Neural Networks, Seattle, WA, USA, 1991, pp. 849-854 vol.1, doi: 10.1109/IJCNN.1991.155290.

[36] MICHAEL J. HUDAK (1992) RCE CLASSIFIERS: THEORY AND PRACTICE, Cybernetics and Systems, 23:5, 483-515, DOI: 10.1080/01969729208927478



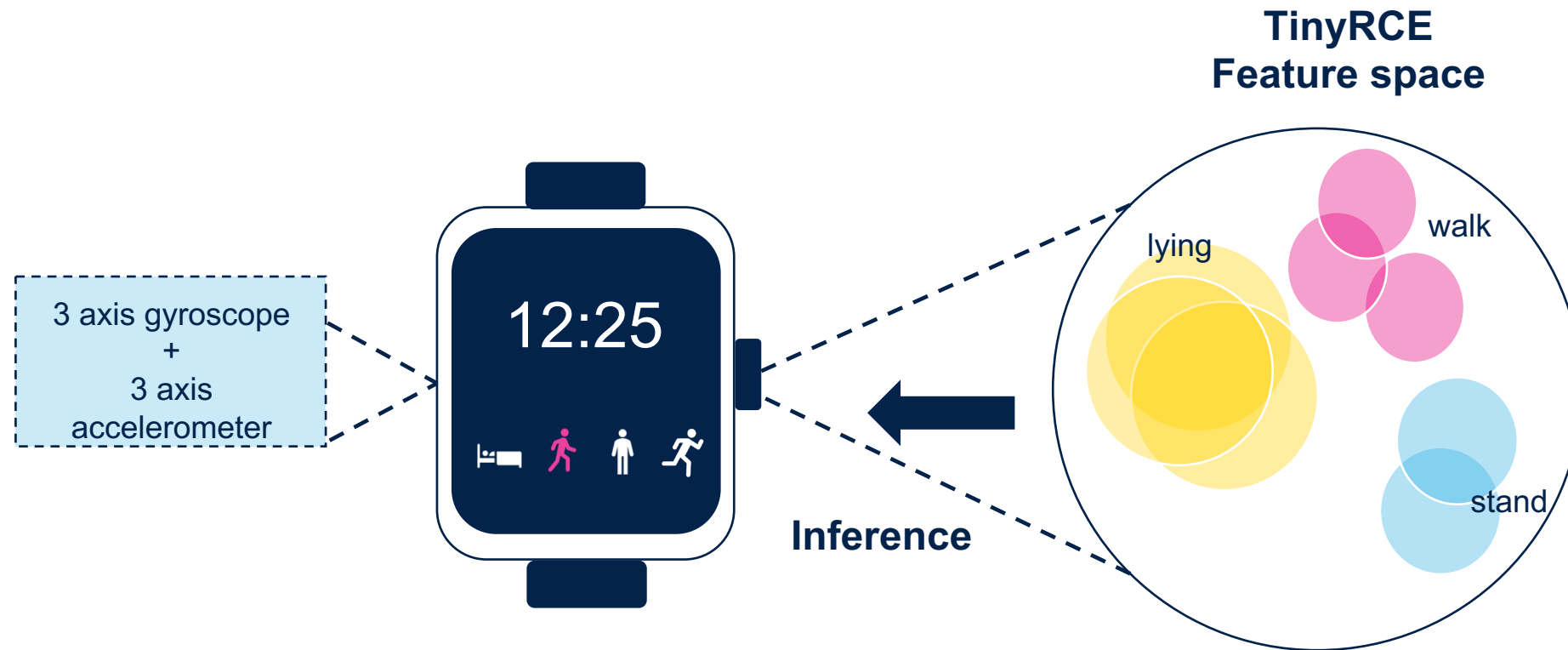
# User annotation process

TinyRCE was meant to detect unknown input patterns, thus triggering the user to provide the corresponding label.



# User annotation process (e.g. HAR)

TinyRCE was meant to detect unknown input patterns, thus triggering the user to provide the corresponding label.



# TinyRCE vs baseline RCE



Learn within  
memory caps

The maximum number of hidden neurons the MCU shall store is depending on the budgeted memory.



Aging

Increased/decreased depending on  $\frac{dist(h_j, x_t)}{R_j}$



Pruning

Prevent uncontrolled instantiation of hidden neurons (causing overflowing available RAM) during the learning phase

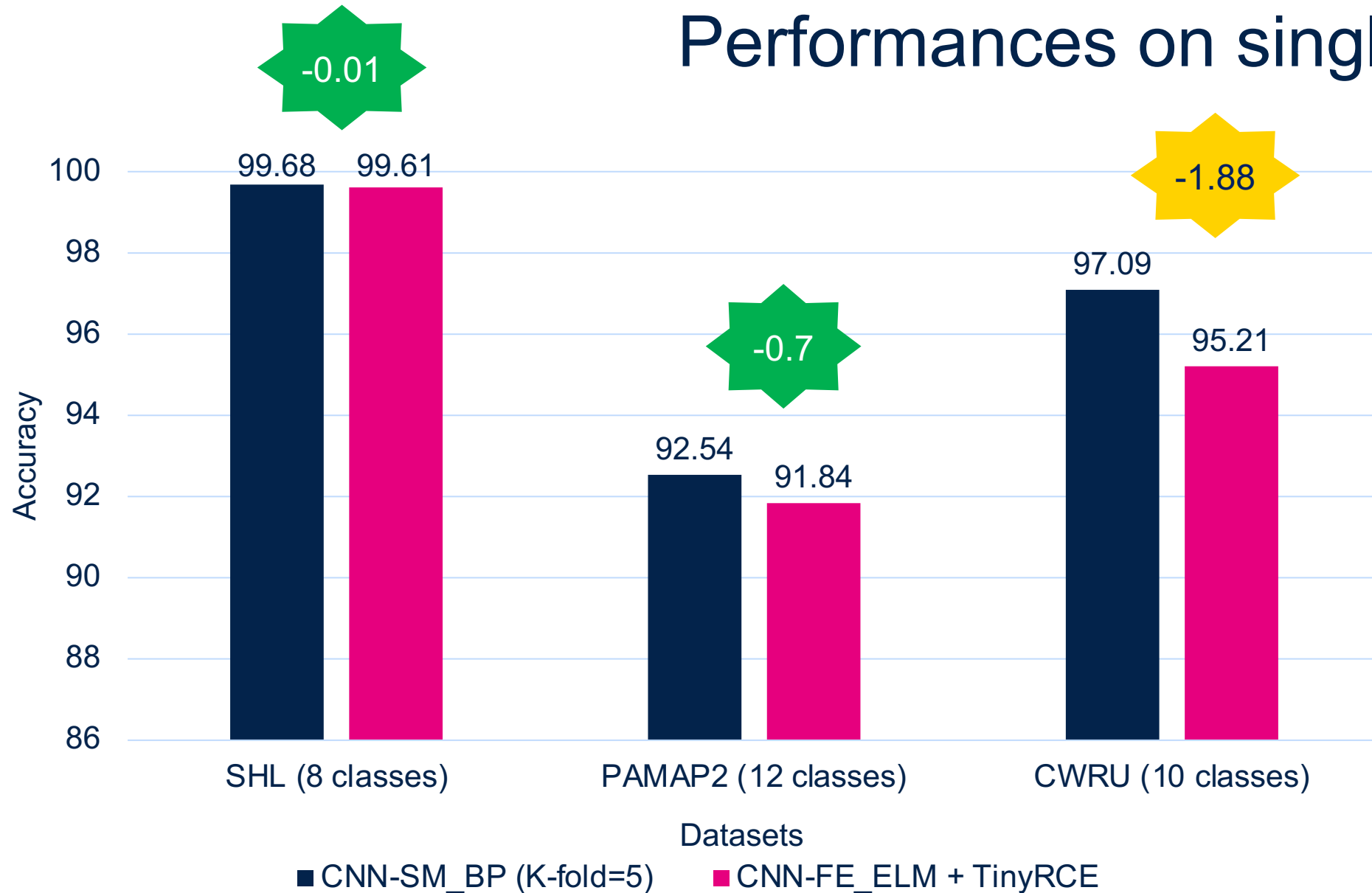


Unambiguous  
assignment

When input falls into multiple hidden neurons regions of influence, the output of the one to which the input falls closer (Hamming or Euclidian distance) is considered



# Performances on single user



# Deployability on MCU

\* random weights

Metrics	CNN-SM_BP (single inference) <u>FP32</u>	CNN-SM_BP (learning) <u>FP32</u> (K fold=1)	CNN-FE_ELM + TinyRCE (single inference) <u>FP32</u>	CNN- FE_ELM + TinyRCE (learning) <u>FP32</u>
MACC	1.213 M	<b>8 G</b>	1.251 M	<b>23.4 M</b>
FLASH (KiB)	76.45	-	76.45 (*)	
RAM (KiB)	13.59	<b>954.57</b>	40.2	
Latency STM32L4 (ms) @80MHz	123.4	<b>813.25 sec</b>	127.2	<b>2.38 sec</b>
Latency STM32H7 (ms) @480 MHz	11.21	<b>73.86 sec</b>	11.55	<b>216</b>

$$MACC\_Inference = h * [(n * 5) + 10]$$

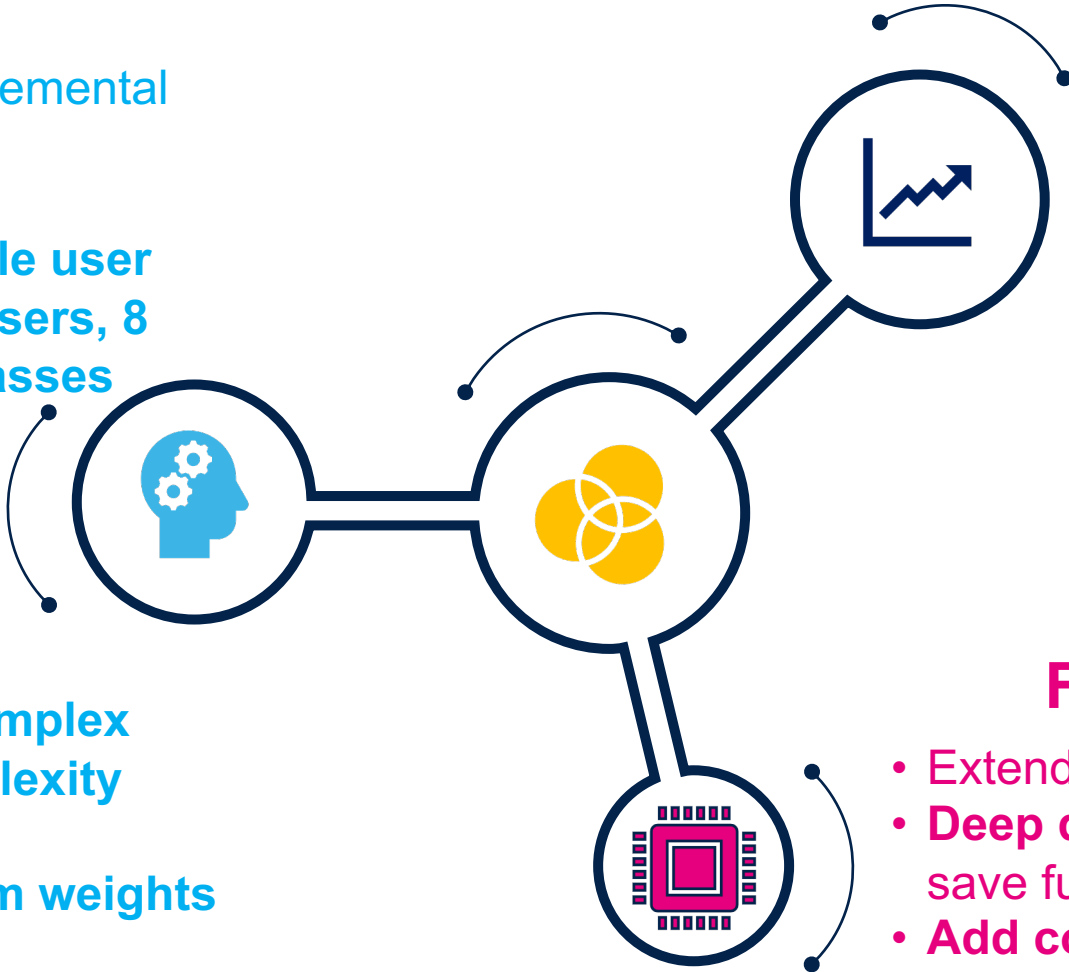
$$MACC\_Learning = (h * [(n * 5) + 10]) * (F * E)$$



# Conclusions

## TinyRCE

- Capable of personalized, incremental learning
- w.r.t. backpropagation
  - 0.01% to -1.88% single user
  - 2% to -3.52% avg 8 users, 8 commons classes
- Deployable on MCUs, w.r.t. backpropagation
  - learning 342x less complex
  - inference same complexity
- Memory
  - same FLASH with random weights
  - 24x less RAM



## Tiny learning

- Open new path to applications
- Save dramatically complexity on the cloud
- Scale faster!
- More reactive Federated Learning

## Future works

- Extend to **MLCommons/Tiny** tasks
- **Deep quantization** of the CNN\_FE to save further memory
- **Add concept drift detection** to trigger learning phase autonomously





life.augmented

**Thanks for listening.**

**Q&A**

**[daniло.pau@st.com](mailto:daniло.pau@st.com)**

# Copyright Notice

This presentation in this publication was presented at the tinyML<sup>®</sup> Research Symposium (March 27,2023). 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)**