

DAC Young Fellows

Accelerating a Visual Wake Word Application on a CGRA

Melina Soysal, Tulika Mitra, Technical University of Munich / National University of Singapore

1 Introduction and Motivation

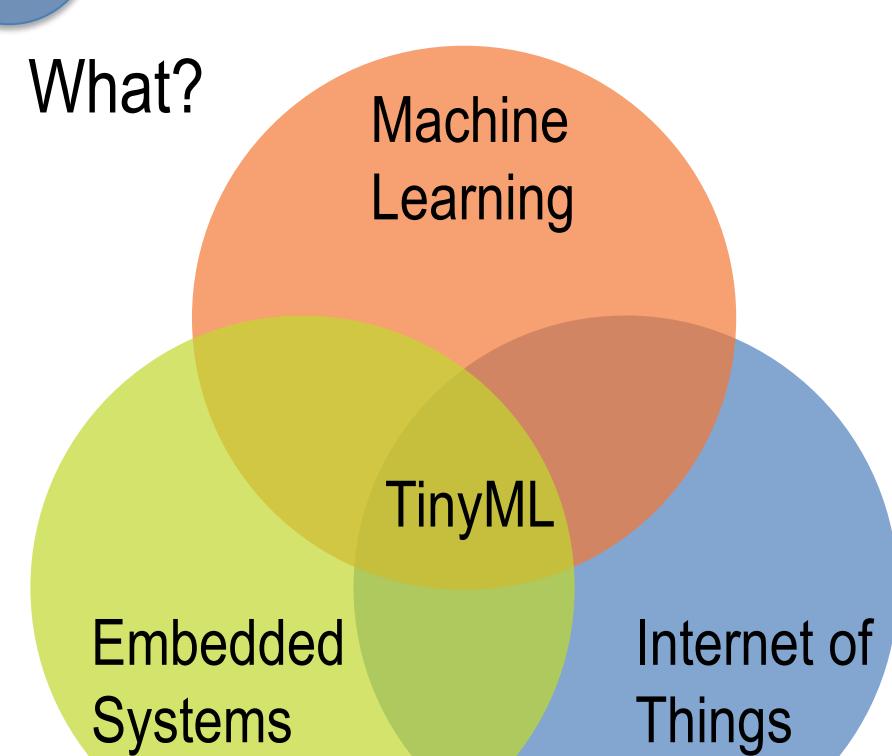
Traditional Machine Learning (ML)

- computationally expensive
- requires a lot of memory
- high energy consumption

Benefits of Microcontrollers (μ Cs)

- cheap
- prevalent
- low-power
- Combining the benefits of a μ C with the power of ML will enable intelligent edge devices and eliminate the necessity for raw data transmission
- Further accelerate this process through task-specific hardware

2 Tiny Machine Learning (TinyML)



What?



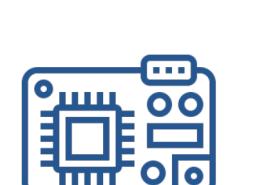
- Machine Learning
- reduce energy consumption
- reduce latency
- increase security and privacy

Why?



- paradigm shift for ML inference
- compute centric in the cloud
- data centric on task-specific μ Cs

Enabled through



Task-specific HW-acceleration



Lightweight ML (reducing redundancy & sparsity)

- TinyML could enable the broad and cost-effective application of machine learning

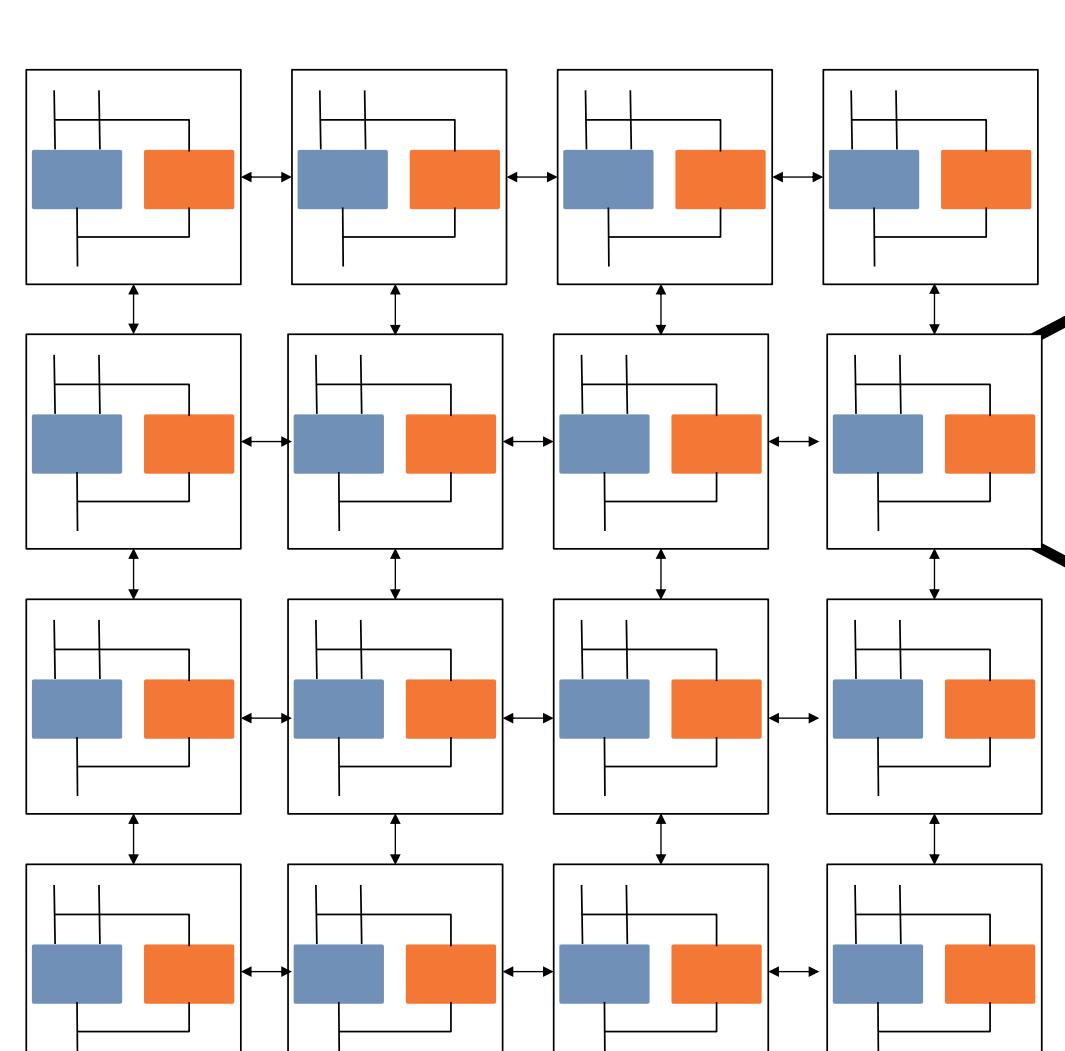
3 Coarse-Grained Reconfigurable Array (CGRA)

reconfigurable system on word level

trade-off between performance & flexibility

well suited for acceleration of loop structures

Structure



- interesting for ML acceleration as computational load often comes from convolution layers which can be represented as loops

4 Visual Wake Word (VWW) application

Idea:



Why Wake Word?

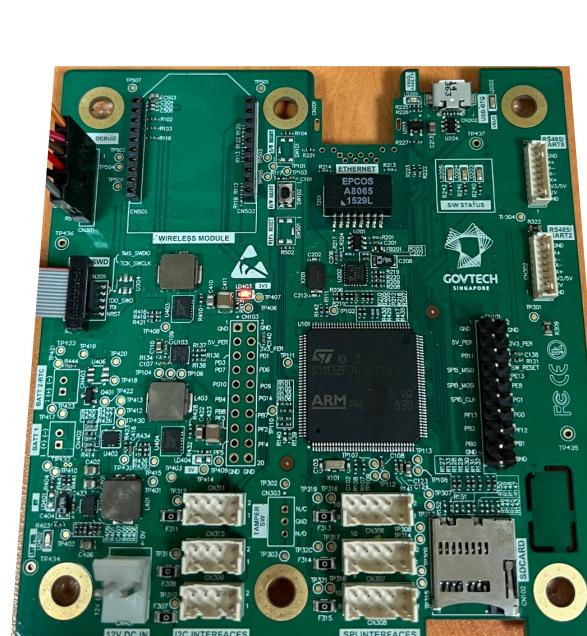
- based on the idea of Wake Words (e.g.: "Hey Siri", "Hey Google", etc.)
- small device detects wake word
- larger, more energy consuming device gets "waken" up

NO person detected

person detected

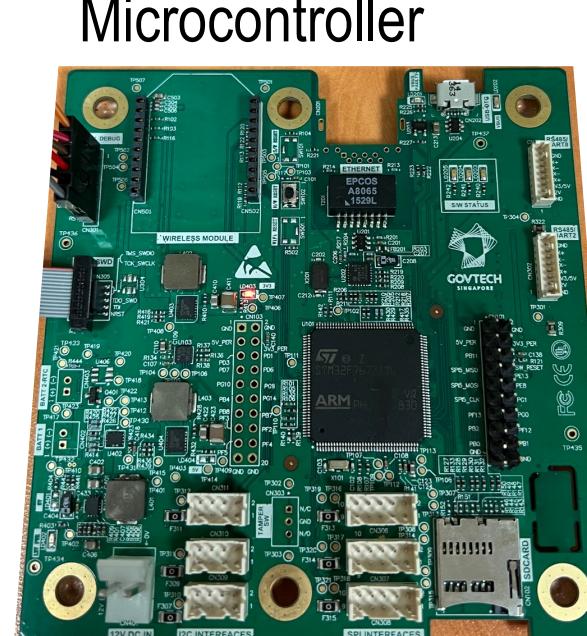
5 Experimental Setup

Current Setup:



DECADA Embedded
MCU: STM32F767ZI based on ARM Cortex M7 CPU
Flash: 2 MB
SRAM: 512 kB

Future Setup:



Microcontroller
connected via QSPI
Accelerator

complete application runs on DECADA Embedded

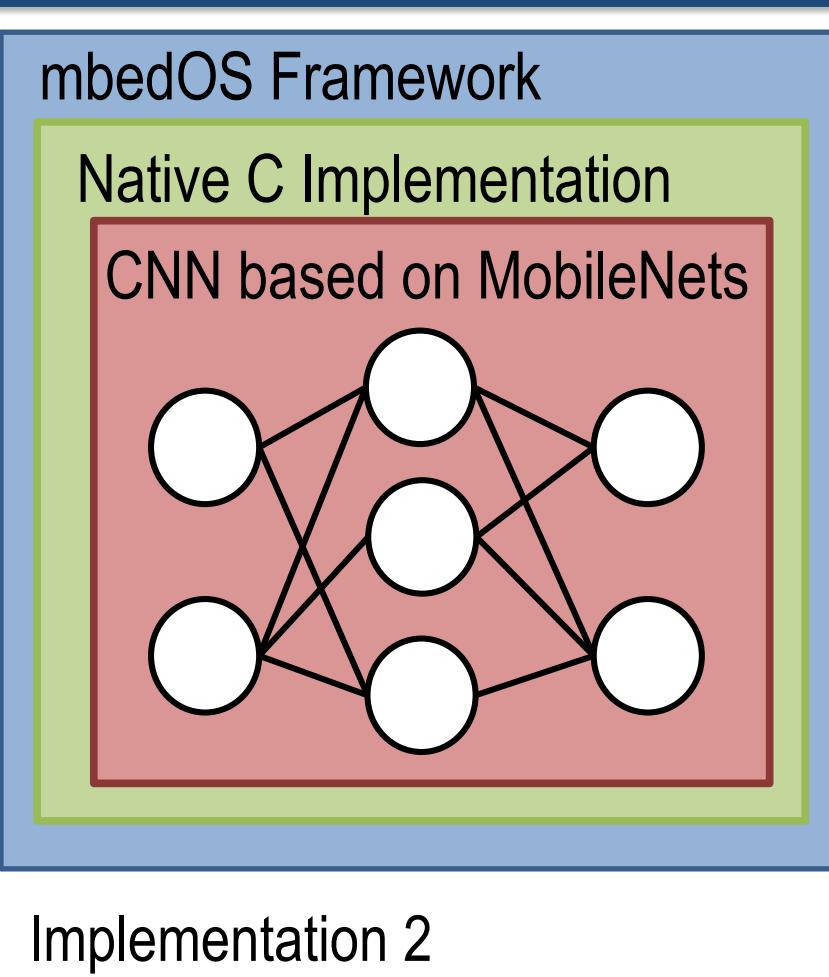
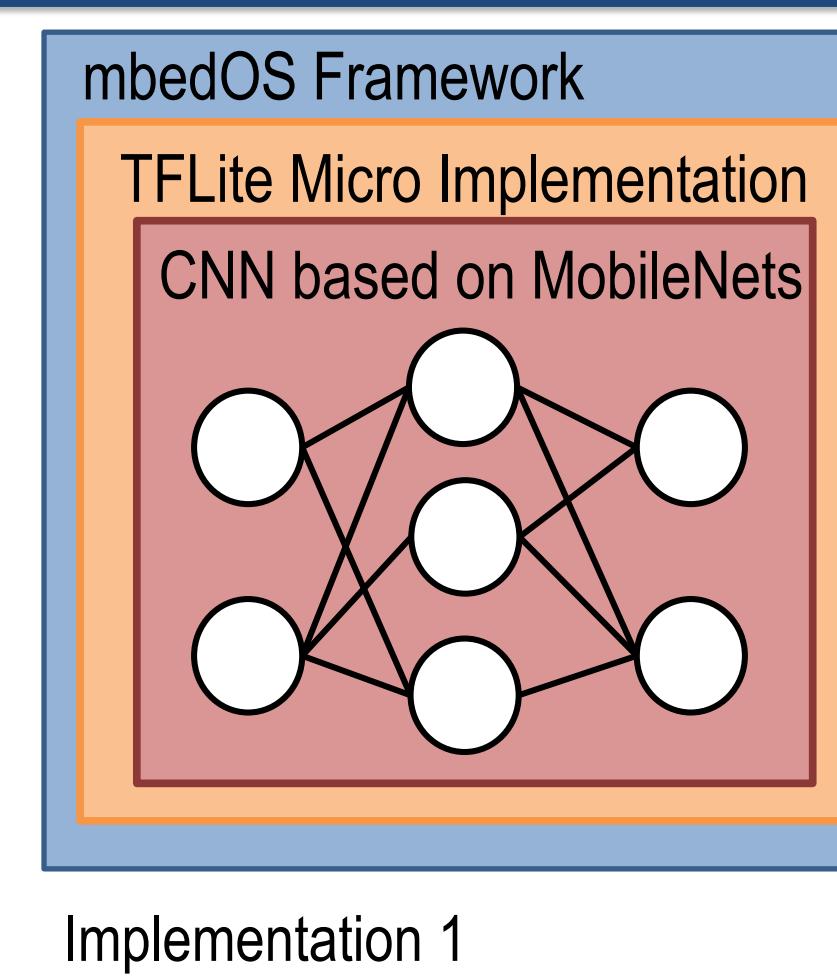
Convolution layers get accelerated on the Hycube¹ CGRA

6 ML Implementation

Person detection is performed by a CNN based on MobileNets²

Depth-wise separable convolutions, GEMM and quantization enable a more lightweight model

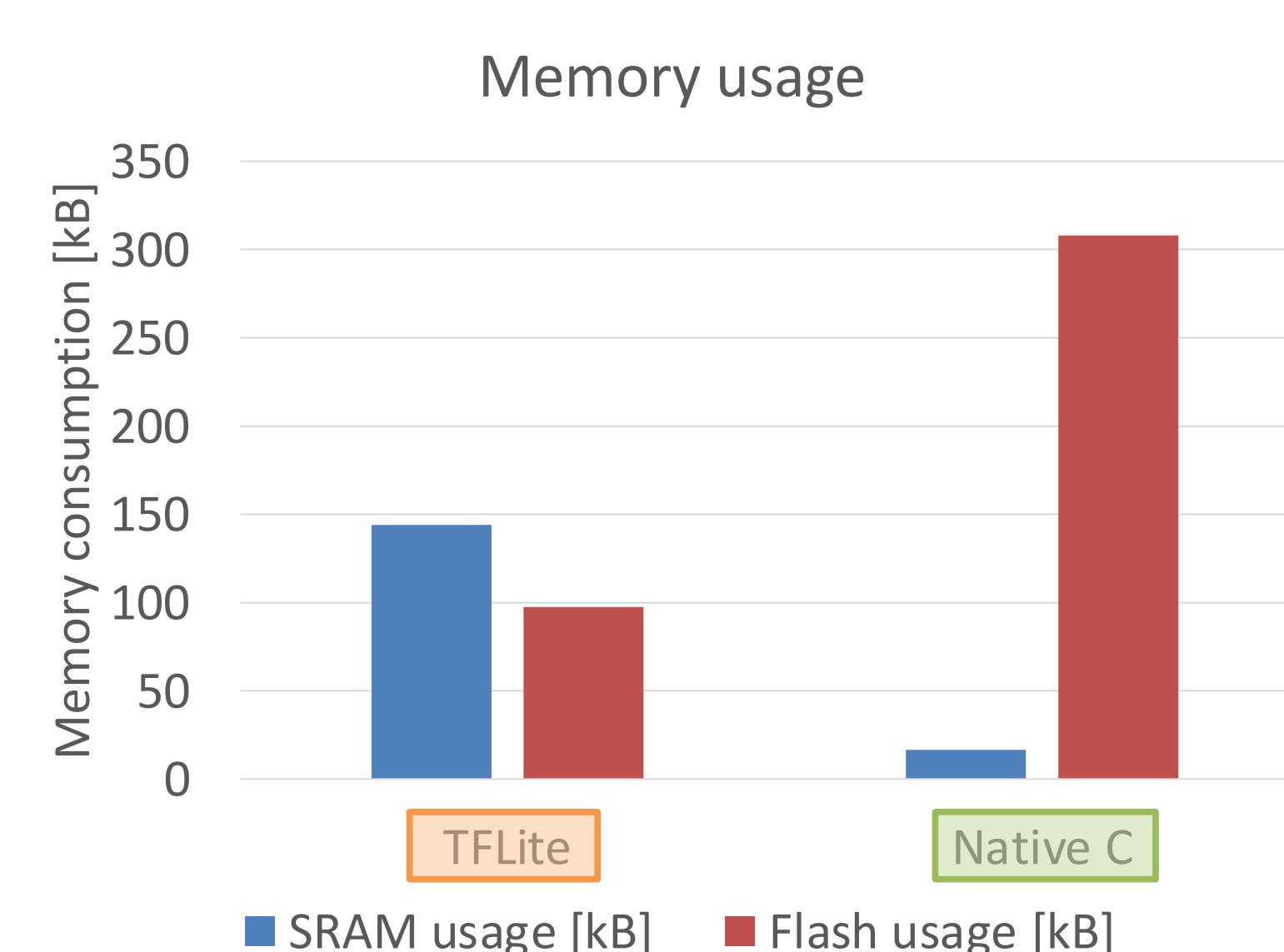
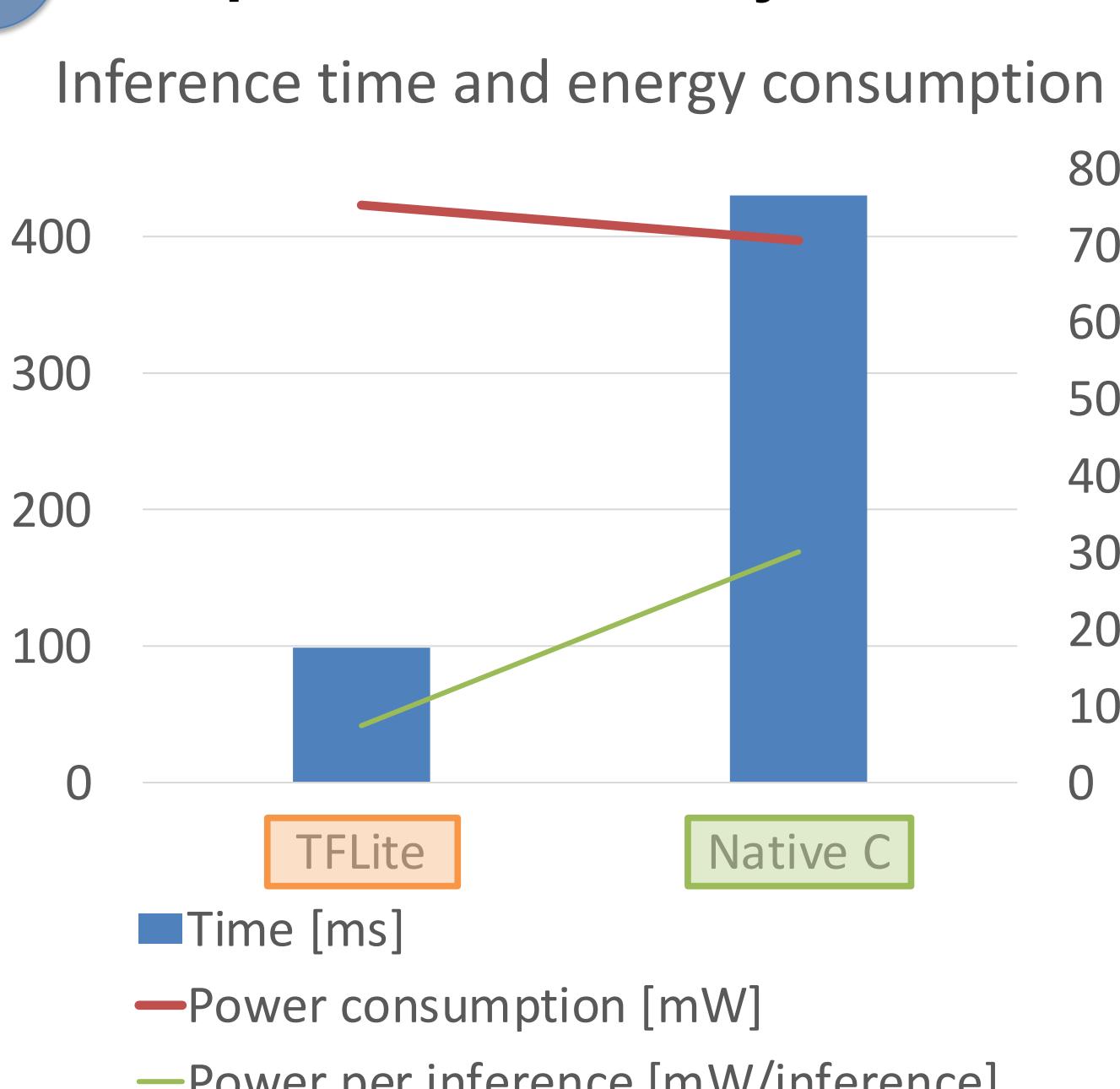
- Two implementations:
 - using the TFLite Micro Framework
 - using only a Native C implementation



Implementation 1

Implementation 2

7 Comparison and analysis of both implementations



The TFLite implementation is faster, more energy efficient per inference and consumes less flash memory.

The Native C implementation will be accelerated with the CGRA, TFLite will serve as baseline.

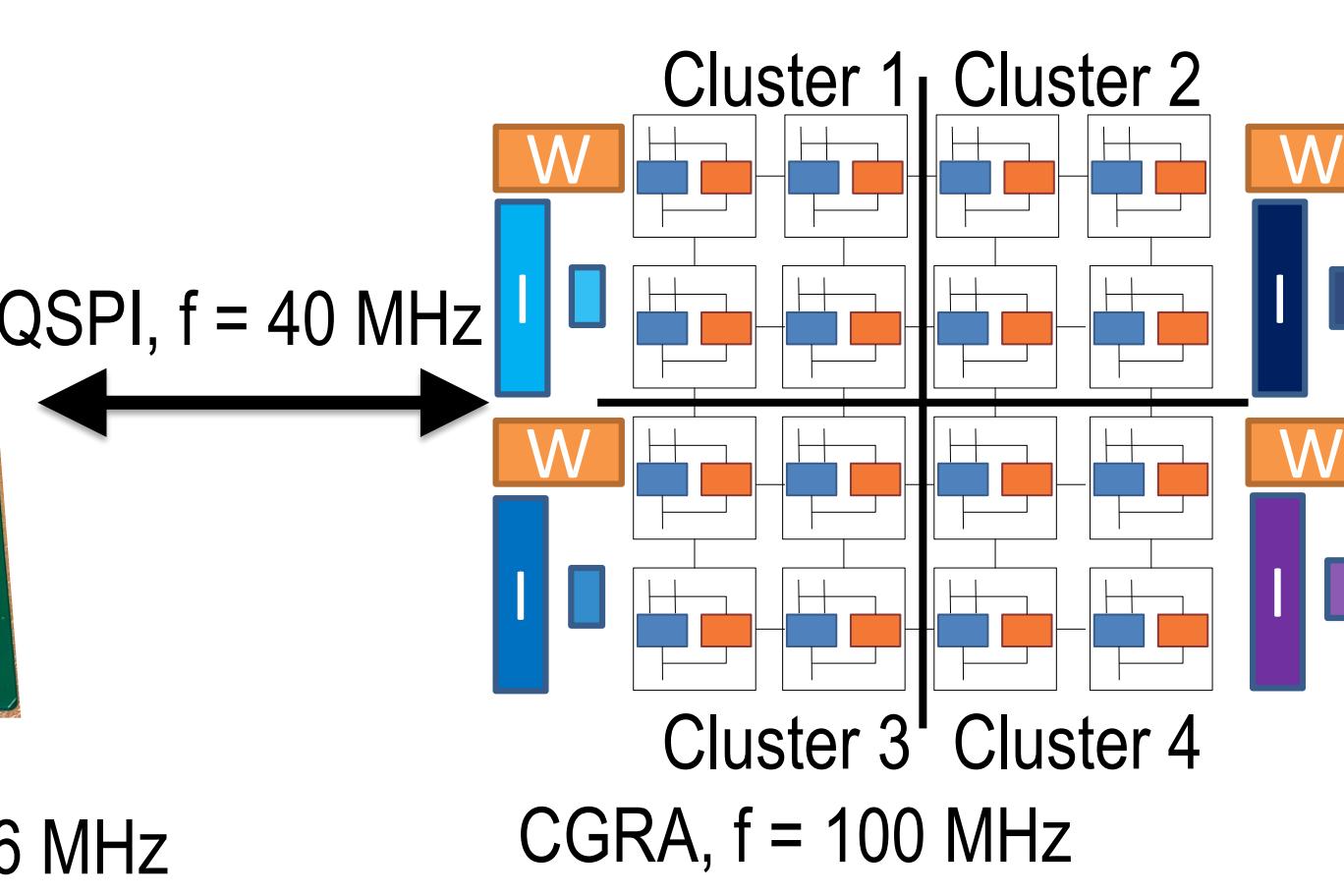
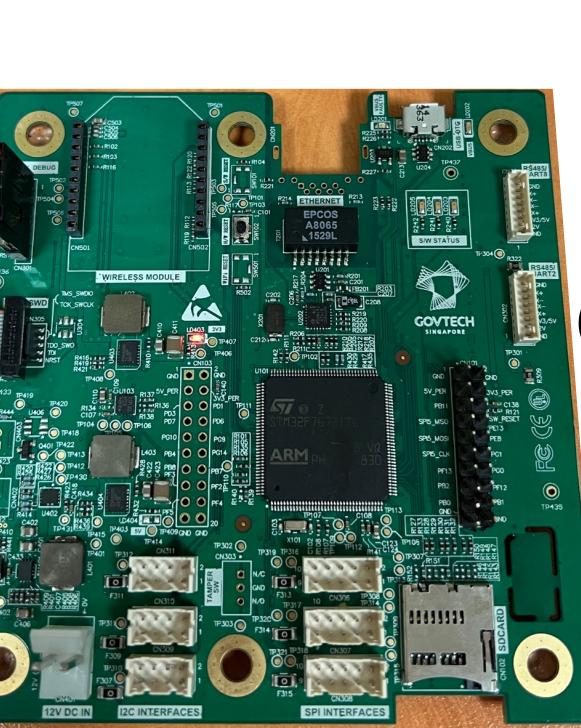
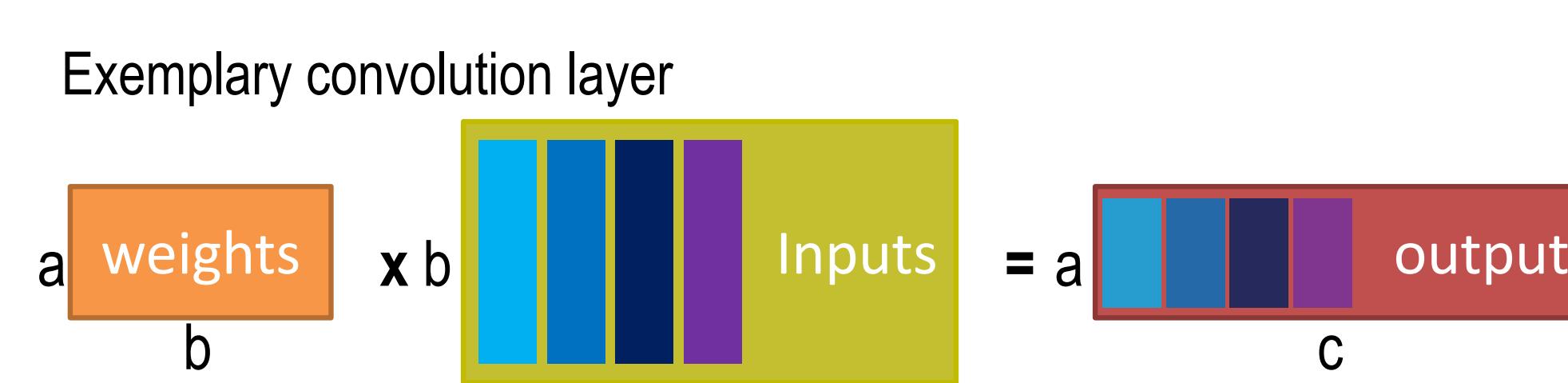
8 Future Work and Conclusion

Acceleration with HyCube¹ CGRA

weights get communicated all at once, input gets communicated in sections due to memory restrictions

for the VWW application some weights are too large to be stored on the CGRA as a whole

- exploding communication overhead (= bottleneck)
- CGRA needs more on-chip memory



9 References

- M. Karunaratne, A. K. Mohite, T. Mitra, and L.-S. Peh. Hycube: A CGRA with Reconfigurable Single-cycle Multi-hop Interconnect, 2017, In Proceedings of the 54th Annual Design Automation Conference 2017
- A.G. Howard et. al., 2017, MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications