



Member Brief: SEI, AI Division, Advanced Computing Lab

Overview of Brief

- **Topic:** Optimization of ML Algorithms for Constrained SWAP Edge Computing in support of Decision Advantage
- **Purpose:** Seeking Mission Partners
- **Description:** Introduce the PHITE project, FY22/23 \$1.5M internally funded applied R&D
- **Host Organization:** Software Engineering Institute (DOD FFRDC at Carnegie Mellon University)
- **Briefer:** Dr. Scott McMillan, Principal Research Engineer (smcmillan@sei.cmu.edu)

Summary of Brief

Discussion of the PHITE (**P**ortable **H**igh-performance **I**nfERENCE at the **T**actical **E**dge) project

- Motivation
- High-level technical approach
- Campus collaboration
- Key Dates & Milestones

Objectives, Tasks, and Deliverables

- AIA SWG Objectives/Tasks
 - #3: Coordinate/Maintain: DoD AI Portfolio
 - Greater understanding of R&D underway with opportunity to leverage existing investment
 - #4: AI Guidance
 - Potential opportunity to inform AIA roadmap through lessons learned and ability to exercise use cases
 - #5: Validate through Partnerships
 - Serves as a Pathfinder for JADC2

Working Group Discussion

- What groups/missions would benefit from PHITE technology? E.g., force protection, ISR, humanitarian operations, predictive maintenance, autonomous platforms.
- What ML algorithms and hardware platforms do these groups/missions deploy?



Copyright 2022 Carnegie Mellon University.

This material is based upon work funded and supported by the Department of Defense under Contract No. FA8702-15-D-0002 with Carnegie Mellon University for the operation of the Software Engineering Institute, a federally funded research and development center.

The view, opinions, and/or findings contained in this material are those of the author(s) and should not be construed as an official Government position, policy, or decision, unless designated by other documentation.

References herein to any specific commercial product, process, or service by trade name, trade mark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by Carnegie Mellon University or its Software Engineering Institute.

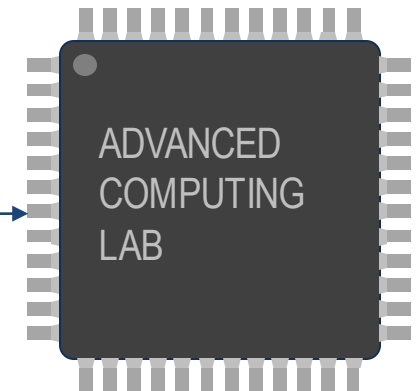
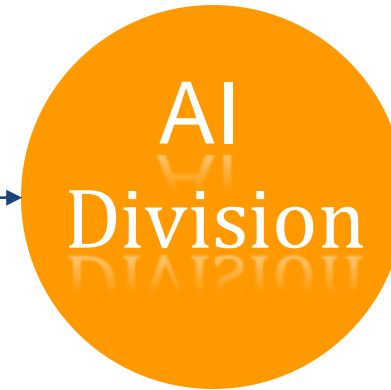
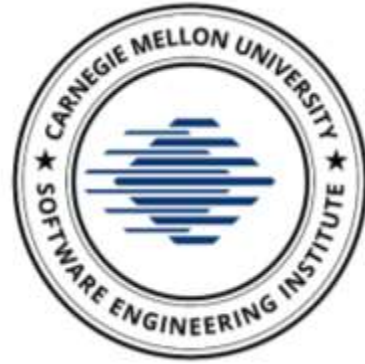
NO WARRANTY. THIS CARNEGIE MELLON UNIVERSITY AND SOFTWARE ENGINEERING INSTITUTE MATERIAL IS FURNISHED ON AN "AS-IS" BASIS. CARNEGIE MELLON UNIVERSITY MAKES NO WARRANTIES OF ANY KIND, EITHER EXPRESSED OR IMPLIED, AS TO ANY MATTER INCLUDING, BUT NOT LIMITED TO, WARRANTY OF FITNESS FOR PURPOSE OR MERCHANTABILITY, EXCLUSIVITY, OR RESULTS OBTAINED FROM USE OF THE MATERIAL. CARNEGIE MELLON UNIVERSITY DOES NOT MAKE ANY WARRANTY OF ANY KIND WITH RESPECT TO FREEDOM FROM PATENT, TRADEMARK, OR COPYRIGHT INFRINGEMENT.

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution. Please see Copyright notice for non-US Government use and distribution.

This material may be reproduced in its entirety, without modification, and freely distributed in written or electronic form without requesting formal permission. Permission is required for any other use. Requests for permission should be directed to the Software Engineering Institute at permission@sei.cmu.edu.

DM22-0230

About us



Vision

To provide the warfighter artificial intelligence solutions that use the latest and most advanced computing technology at all scales in a timely, reliable, and safe and secure manner that maximizes warfighter effectiveness.

Mission

Apply our world class expertise to solve advanced computing problems in artificial intelligence at all scales – edge, cloud, and HPC – that enable continuous improvement of the warfighter's ability to defend our nation.

AI at the Tactical Edge

PHITE: Portable High-performance Inference at the Tactical Edge

- Wearable technology
- Autonomous UAVs/UGVs for
 - ISR
 - HADR
- Ground/Unattended Sensors
- Smart/Cognitive Radio
- TAK Display

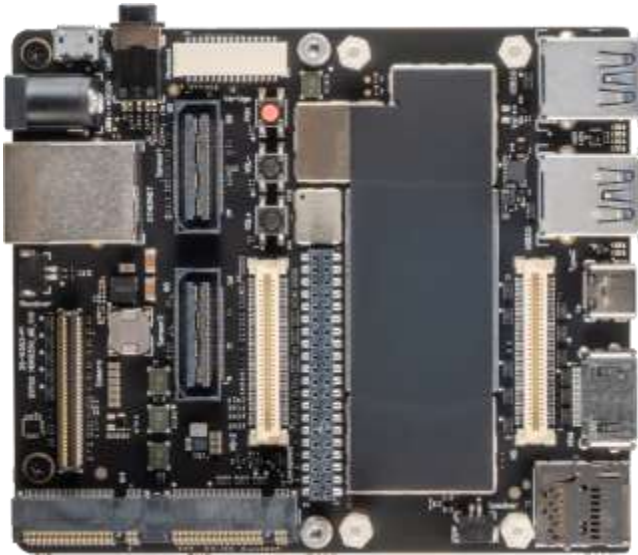
Key challenges

- SWAP
- Battery limitations
- Mission time



PHITE Goal: do more with less.

Hardware at-a-glance



System:	Qualcomm Snapdragon 888 Development Kit	Raspberry Pi Model 4-B SOC	Raspberry Pi PICO Microcontroller	Your hardware here...
Power	3-8W platform, <TBD>W Hexagon	~3W (idle) – 7W (4 cores)	6mW - 330mW	
Cost	\$1,349 (for a dev board)	\$35-\$75	\$4	
Memory	<TBD>GB dedicated to Hexagon	2/4/8 GB SDRAM	264KB RAM (2MB flash)	
Processor	Kryo 680 CPU, Adreno 660 GPU, Hexagon 780 AI Accelerator	Broadcom BCM2711, ARM v7I, Cortex-A72 (4 cores)	RP2040: ARM Cortex-M0+ (dual core)	
Peak	26 TOPS (Hexagon)	13.5 gigaFLOPS	266 megaFLOPS	
Model/Size	<TBD>	AlexNet / 60M parameters	MobileNet V2 / 3M parameters	

Background

Much of our efforts are targeted at extending the 2018 research on direct convolutions.

Lower memory footprint (no overhead)

More efficient computation

High Performance Zero-Memory Overhead Direct Convolutions

Jiyuan Zhang¹ Franz Franchetti¹ Tze Meng Low¹

Abstract

The computation of convolution layers in deep neural networks typically rely on high performance routines that trade space for time by using additional memory (either for packing purposes or required as part of the algorithm) to improve performance. The problems with such an approach are two-fold. First, these routines incur additional memory overhead which reduces the overall size of the network that can fit on embedded devices with limited memory capacity. Second, these high performance routines were not optimized for performing convolution, which means that the performance obtained is usually less than conventionally expected. In this paper, we demonstrate that direct convolution, when implemented *correctly*, eliminates all memory overhead, and yields performance that is between 10% to 400% times

Performance normalized to OpenBLAS GEMM on AMD PileDriver
4.0 GHz, 4/4 cores/threads

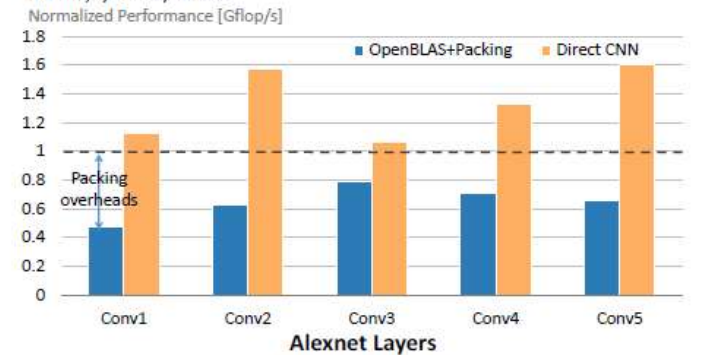


Figure 1. High performance direct convolution implementation achieves higher performance than a high performance matrix multiplication routine, whereas matrix-multiplication based convolution implementations suffers from packing overheads and is limited by the performance of the matrix multiplication routine

In *International Conference on Machine Learning*, pp. 5776-5785. PMLR, 2018.

Key Dates and Milestones

Major Milestones	Fiscal Year	
	22	23
Select and analyze three edge applications (at least one DoD problem)	X	
Select and analyze hardware capability of at least three HW devices	X	
Develop performance models for HW devices	X	
Develop optimized algorithms for targeted HW (at least 1 in FY22)	X	X
API and library development, application development on HW		X
Evaluation of applications against baseline system, benchmarking		X
Demonstration of DoD application with streaming data		X

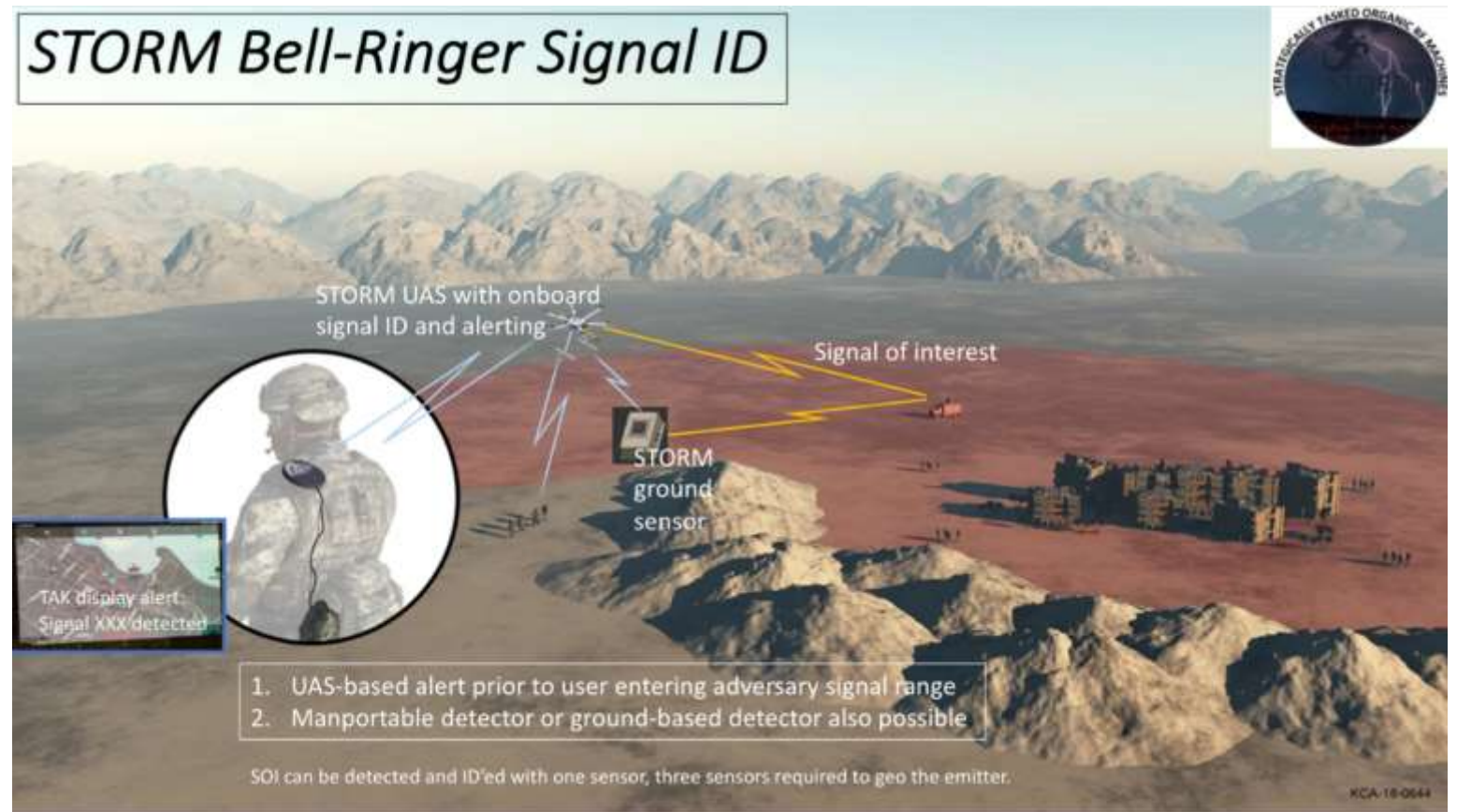
Hardware-Software Co-design for Edge AI Applications

“In general, **compute** can improve mission time and lower energy consumption by as much as 5X.”

Boroujerdian, Behzad, et al. "Why compute matters for UAV energy efficiency?" (2018)

Objective 1: Application dependent SoC co-designs for stakeholders with accelerators that are 2x more energy efficient than GPU solutions and are 30% better in mission metrics.

Objective 2: A co-design flow methodology and **codriver** implementation that includes four tools per pipeline stage and enables co-design of RTL for an application that includes accelerator IP *within two months*.



The 'Ideal' mission partner

- Needs ML algorithms to run on low-power, mobile and embedded hardware devices (use case)
- Can provide access to algorithms and data to validate results
- Can provide access to hardware exemplars or specifications of COTS platforms
- Is engaged, focused and vested in outcome
- Is a fearless champion

Backup

PHITE Software Architecture

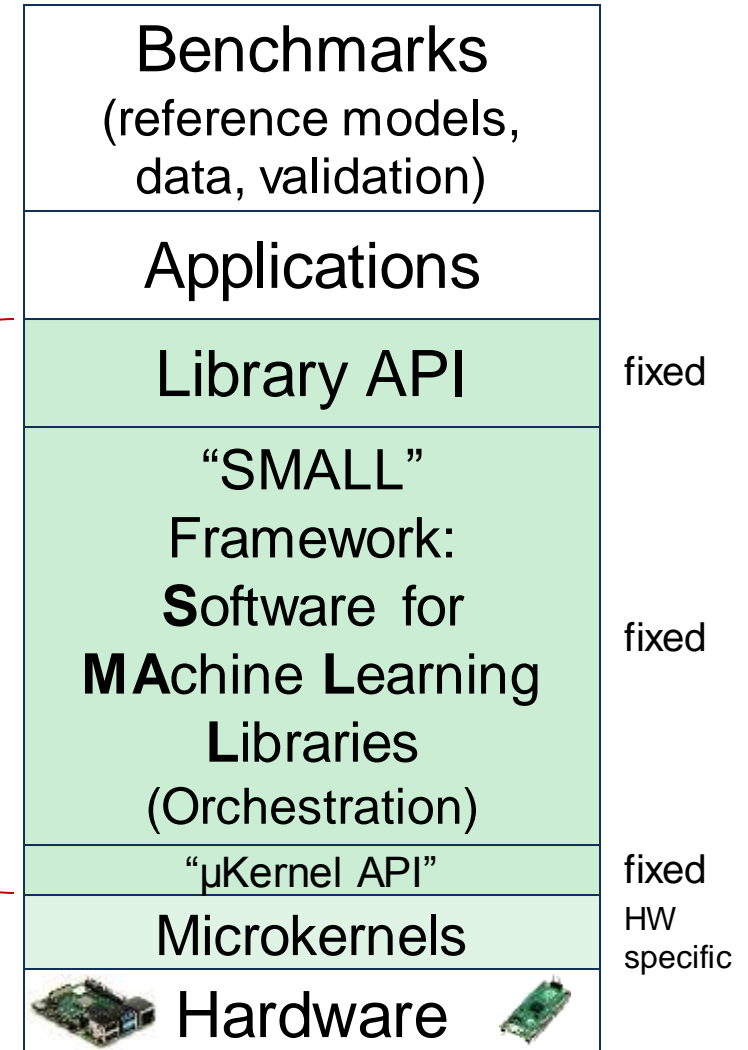
Benchmark and application selection:

- **Open-source:** TinyMLPerf* benchmark applications:
 - image classification, visual wake words, audio wake words, object detection
- **CUI / ITAR:** <Your applications and data here>
 - E.g., target identification, anomaly detection.

Library development (open source): a BLIS[†]-like layered architecture

- **Top-level API:** Tensorflow 2 (as much as possible)
- **Middleware:** SMALL framework (under development)
- **Low-level:**
 - A microkernel API for small set of key machine learning primitives
 - Hardware specific codes, including assembly where necessary

Open source



- <https://github.com/mlcommons/tiny>
- <https://mlcommons.org/en/inference-tiny-05/>

[†] “BLIS: A Framework for Rapidly Instantiating BLAS Functionality”, (TOMS) ACM Transactions on Mathematical Software, vol. 41, no. 3, June 2015, pp 1-33.

Team



Dr. Scott McMillan, CMU/SEI/AI, **PI**

Advanced computing, parallel and distributed algorithms, graph analytics, data-intensive computing, interactive data science at scale.



Prof. Tze Meng Low, CMU/ECE, **PI**

High-performance algorithms using formal methods and analytical models; performance portability through capture of interaction between software algorithms and hardware features; code generation and libraries for emerging domains.



Oren Wright,
CMU/SEI/AI,
CMU/ECE PhD

Applied AI/ML; graph signal processing; moving beyond ad hoc neural networks; all things Bayesian; “how to make AI fail gracefully...not catastrophically.”



Upasana Sridhar,
CMU/ECE PhD

Analytical modelling for graph algorithms and machine learning inference; performance modelling, formal methods.

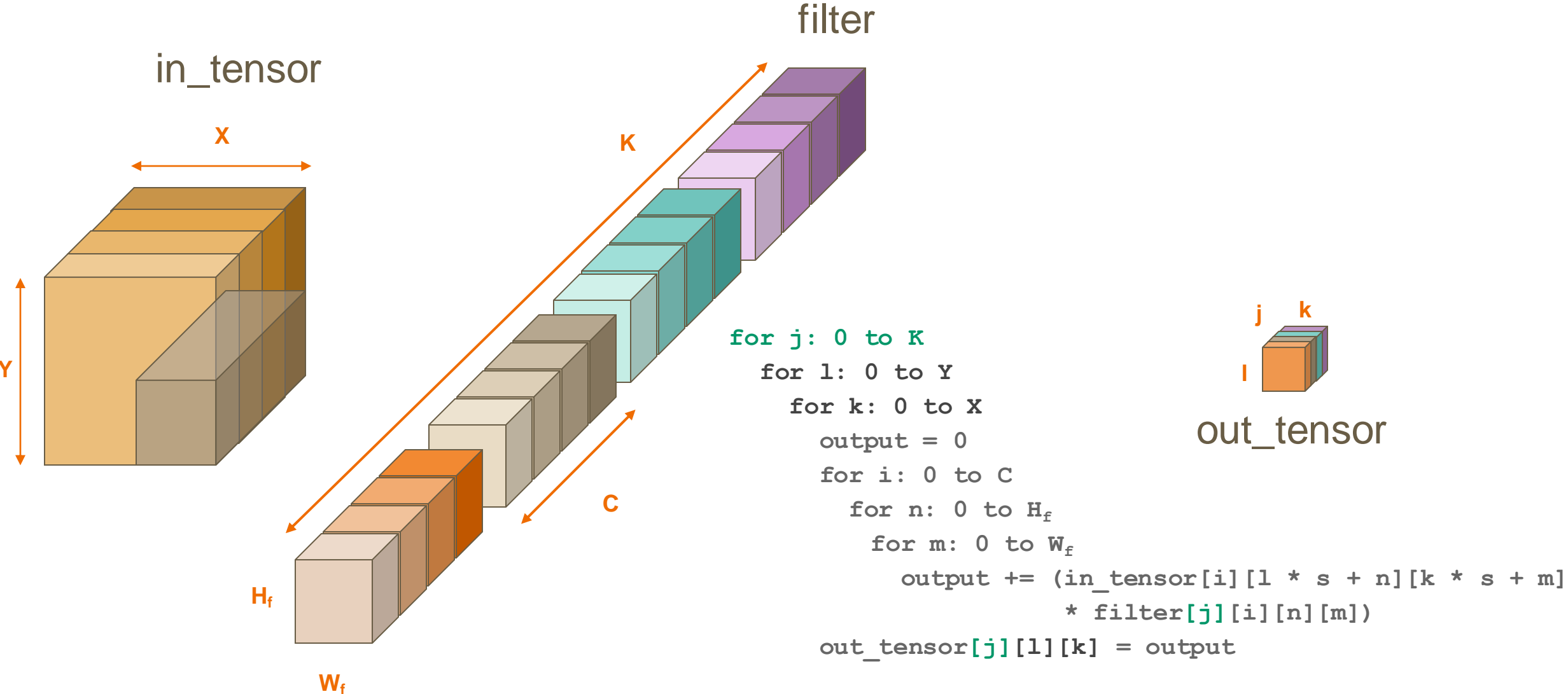


Nicolai Tukanov,
CMU/ECE PhD

High-performance computing, performance modeling, hardware accelerators.

- Navya Chandra – ECE master’s independent study: “Fused convolution on Pi Pico”
- Pankti Rajesh Shah – master’s student started in Q2FY22

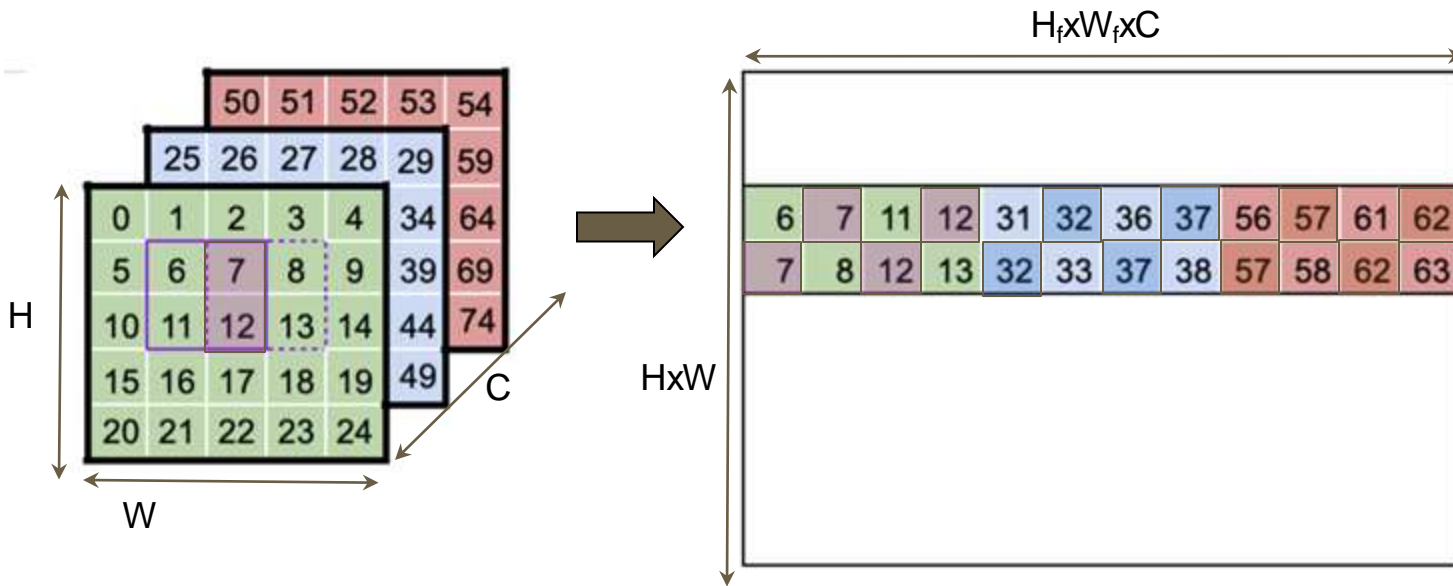
Convolution Operation - End to End



Convolutions as GEMM

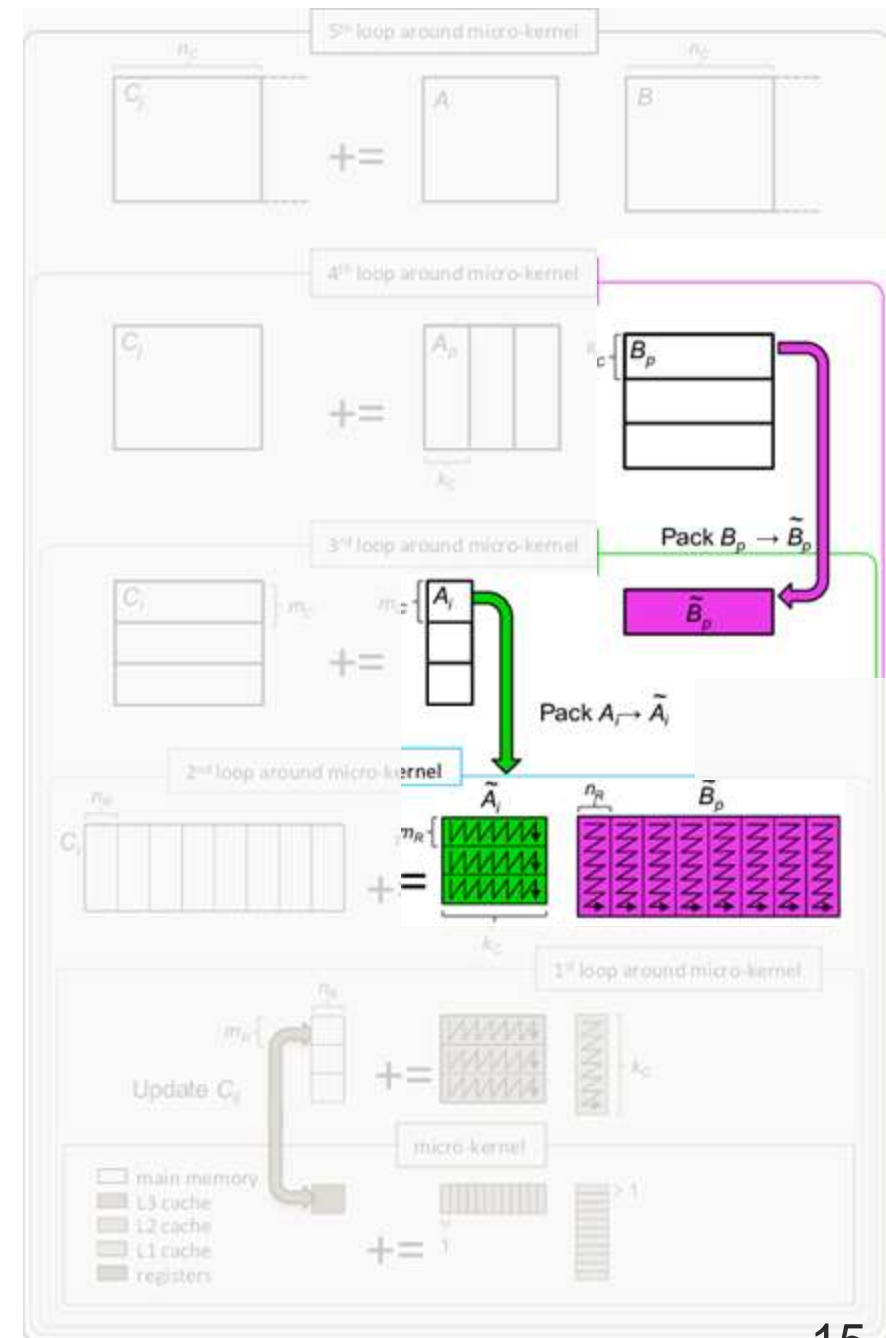
- GEMM: BLAS dense **G**eneral **M**atrix-**M**atrix multiply
- Packing requires additional memory

GEMM packing



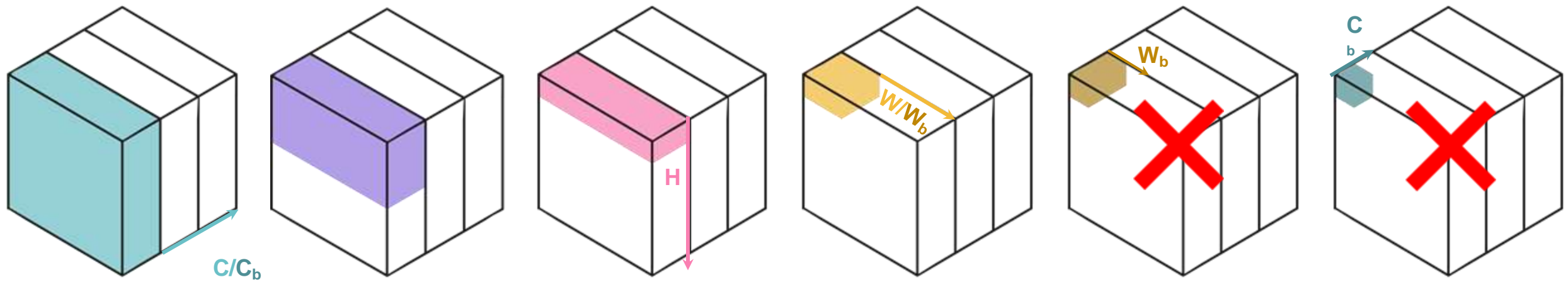
Data Duplication

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution. Please see Copyright notice for non-US Government use and distribution.



Fusing Loop Nests

- Fusion may be possible at different levels



Decreasing (Working Set Size) WSS

Fuse each block

Fuse each row
pool

Fuse each row

Fuse each
pixel block

Fuse each pixel

Fuse each
channel

Highest Level
(fuse outermost loop)

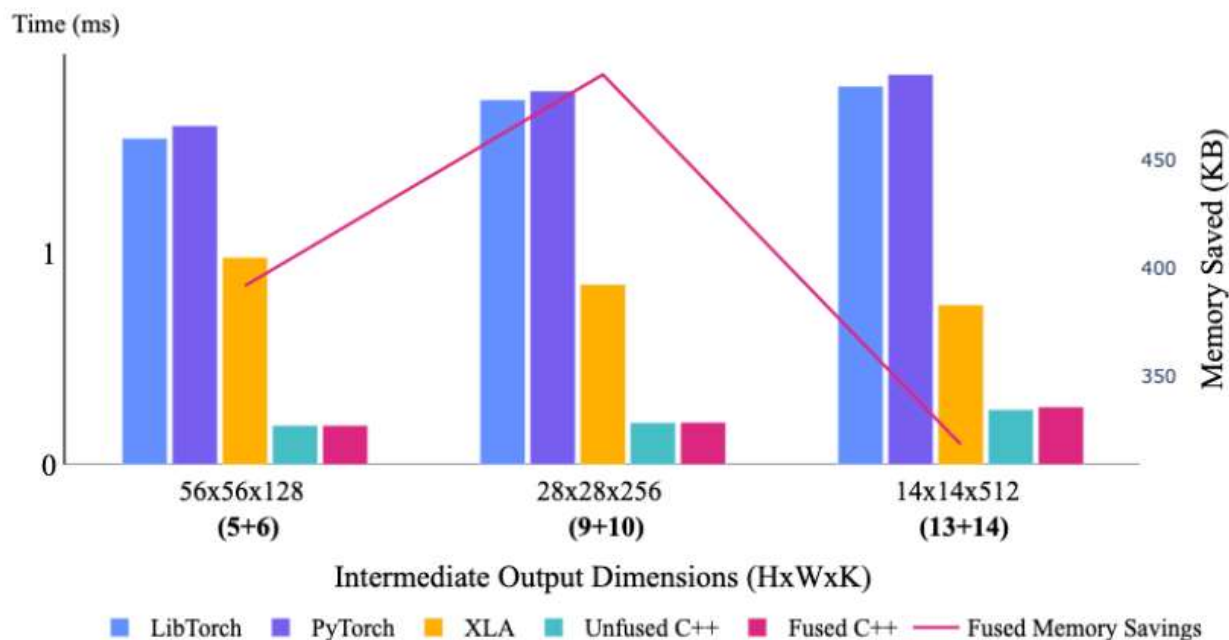
[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution. Please see Copyright notice for non-US Government use and distribution.

Lowest Level
(fuse up to innermost loop)

Fused Loop Structure

- LOTS of knobs to tune: F_C , O_h , O_w , F_h , F_w , I_{cb} , O_{wb} , etc.

MobileNet - 1x1 Convolution Fused with 3x3 Dwise Convolution, stride = 1



Algorithm 2 Fused abstract deep learning layer loop structure

```

1: for  $g \leq G$  do
2:   for  $j \leq K$  do
3:
4:     for  $i \leq F_C$  do ▷ Fusion in Last iteration
5:
6:       for  $k \leq O_h$  do ▷ Transformation Req for Fusion
7:         for  $l \leq O_w$  do
8:           for  $x \leq F_h$  do
9:             for  $y \leq F_w$  do
10:              for  $ii \leq I_{cb}$  do
11:
12:                for  $ll \leq O_{wb}$  do
13:                  for  $jj \leq O_{cb}$  do
14:                    ReductionOp( $O_0, I, F_0$ )
15:                  end for
16:                end for
17:                ▷ Single Element Reduction
18:              for  $ll \leq O_{wb}$  do
19:                for  $jj \leq O_{cb}$  do
20:                  ReductionOp( $O, O_0, F_1$ )
21:                end for
22:              end for
23:            end for
24:          end for
25:        end for
26:        ▷ Channel Reduction
27:      for  $x \leq F_h^1$  do
28:        for  $y \leq F_w^1$  do
29:           $ii = 1$ 
30:          for  $ll \leq O_{wb}$  do
31:            for  $jj \leq O_{cb}$  do
32:              ReductionOp( $O, O_0, F_1$ )
33:            end for
34:          end for
35:        end for
36:      end for
37:    end for
38:  end for
39: end for
40: end for
41: end for=0
  
```