



# **GAP***flow*

Presenter:

**Marco Fariselli**

Embedded ML Engineer

[marco.fariselli@greenwaves-technologies.com](mailto:marco.fariselli@greenwaves-technologies.com)

# GreenWaves Technologies

- Fabless semiconductor startup founded in 2014.
- We design and sell **extreme performance** processors for **energy constrained devices**
- 45 people, HQ in **Grenoble**, France
- Offices in **Bologna**, Italy, **Shanghai**, China, **Copenhagen**, Denmark. Global sales footprint.



Best hardware product  
Embedded World 2023



Embedded Technologies Award 2023  
Les Assises de l'Embarqué

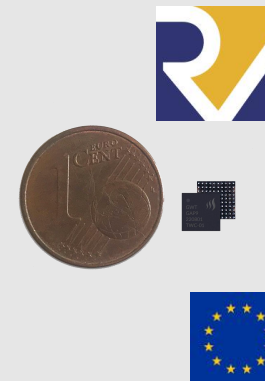
## GAP8

In production since 2020  
one of the very first  
commercially available  
RISC-V processor and AI  
microcontroller



## GAP9

Second generation  
ultra-low power AI and  
DSP enabled  
Microcontroller



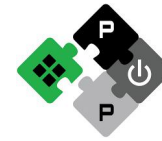
Gartner

COOL  
VENDOR  
2019

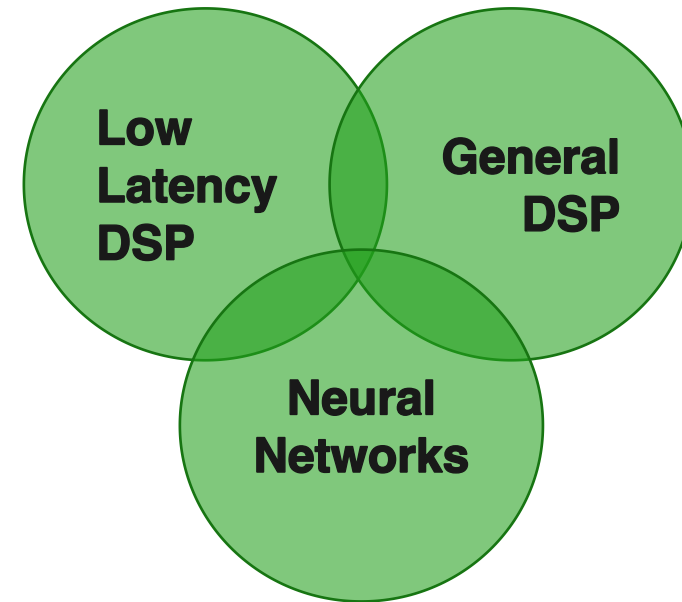
Cool Vendors in AI Semiconductors,  
Alan Priestley, Saniye Alaybeyi, April  
29, 2019.



# Foundations of GAP



- Born from research groups:
  - **PULP** (Unibo & ETH Zurich)
  - **RISC-V**
- Push **Edge AI** to the limits
  - Specialized HW
  - Optimized SW
  - 10-100 GOPS @ <10mW always-on
- **Programmability**
  - RISC-V open source ISA
  - High level tools



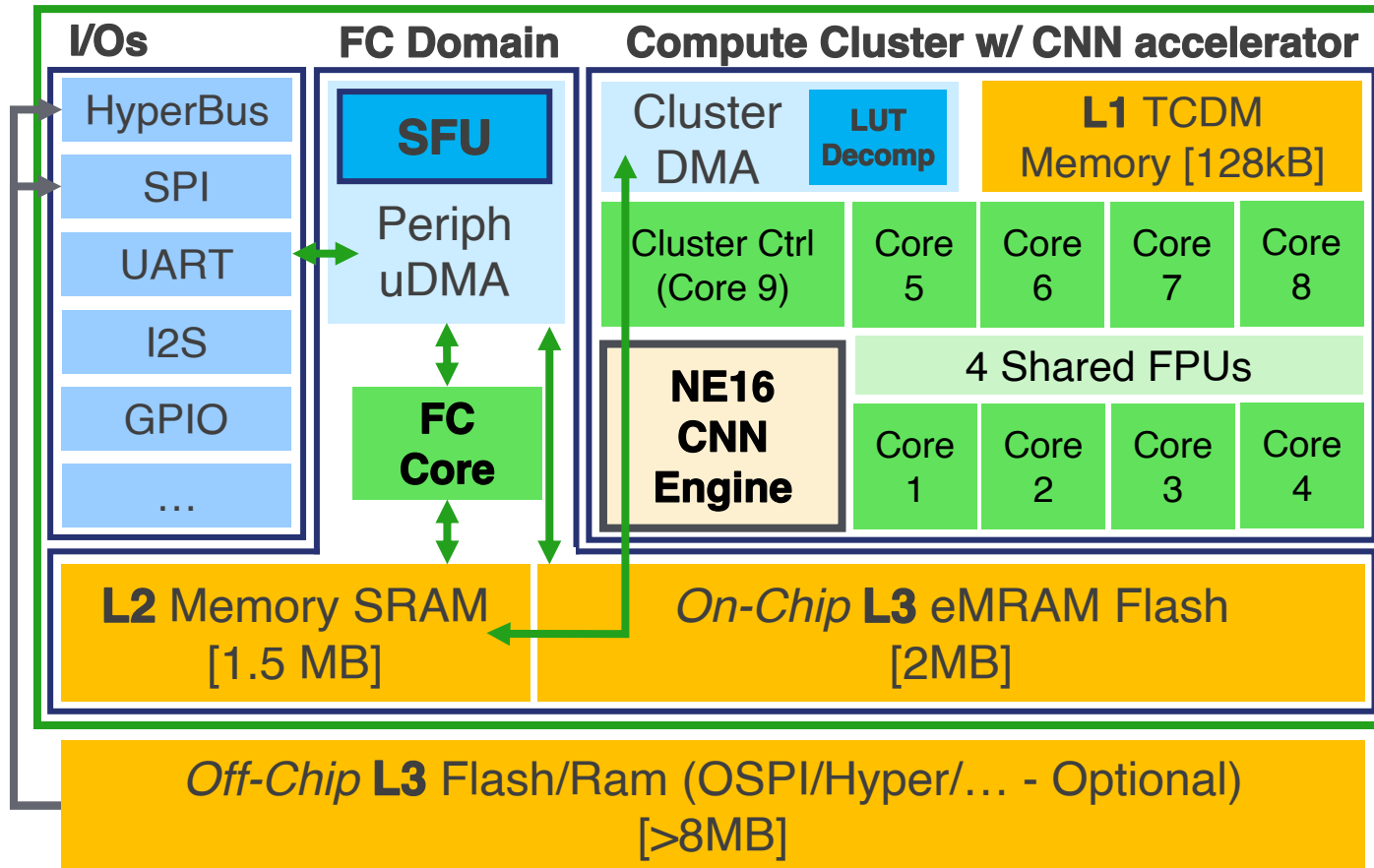
# Outline

- **GAP Architecture**
- **GAPflow**
  - **NNtool**: Graph optimizations / Quantization
  - **Autotiler**: Memory management
- **NE16**
- **Hands-on**

# GAP Architecture

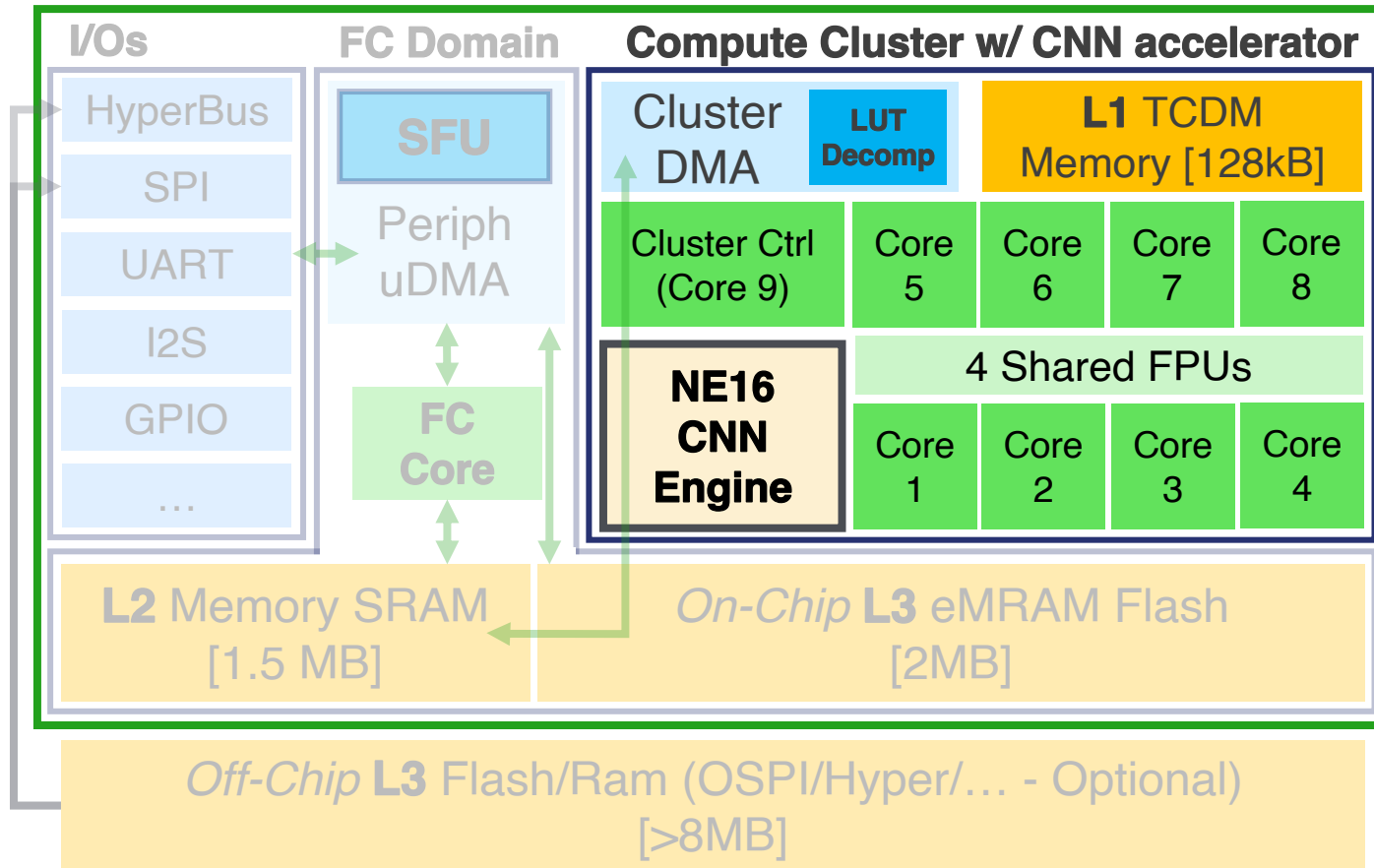
# GAP9 Architecture

## GAP9 SoC



# GAP9 Architecture

## GAP9 SoC



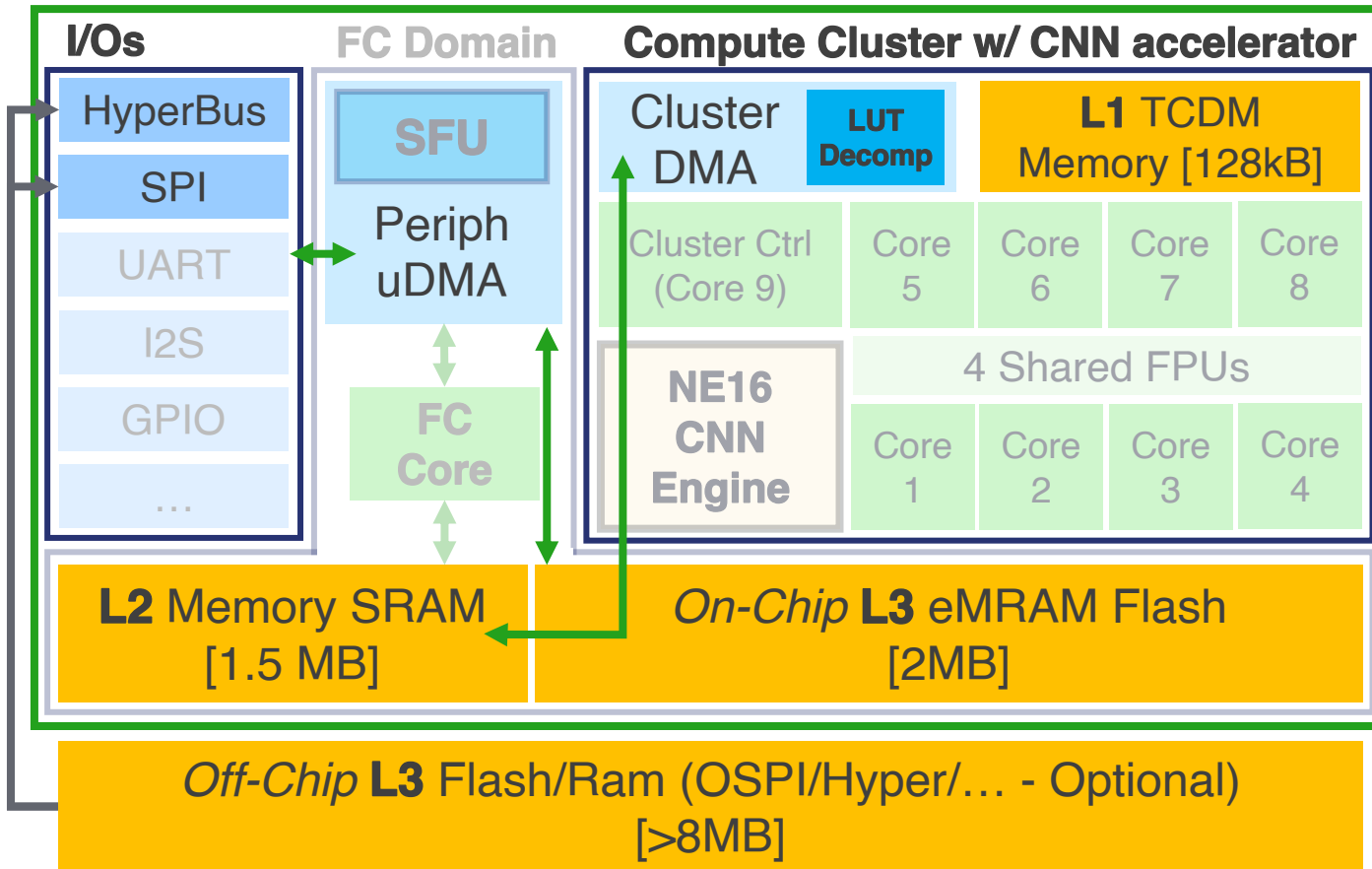
### Hierarchical Compute paradigm

- 4 independent frequency domains: **FC - I/Os - Cluster - SFU**
- “Turn-on when you need”



# GAP9 Architecture

## GAP9 SoC



### Hierarchical Compute paradigm

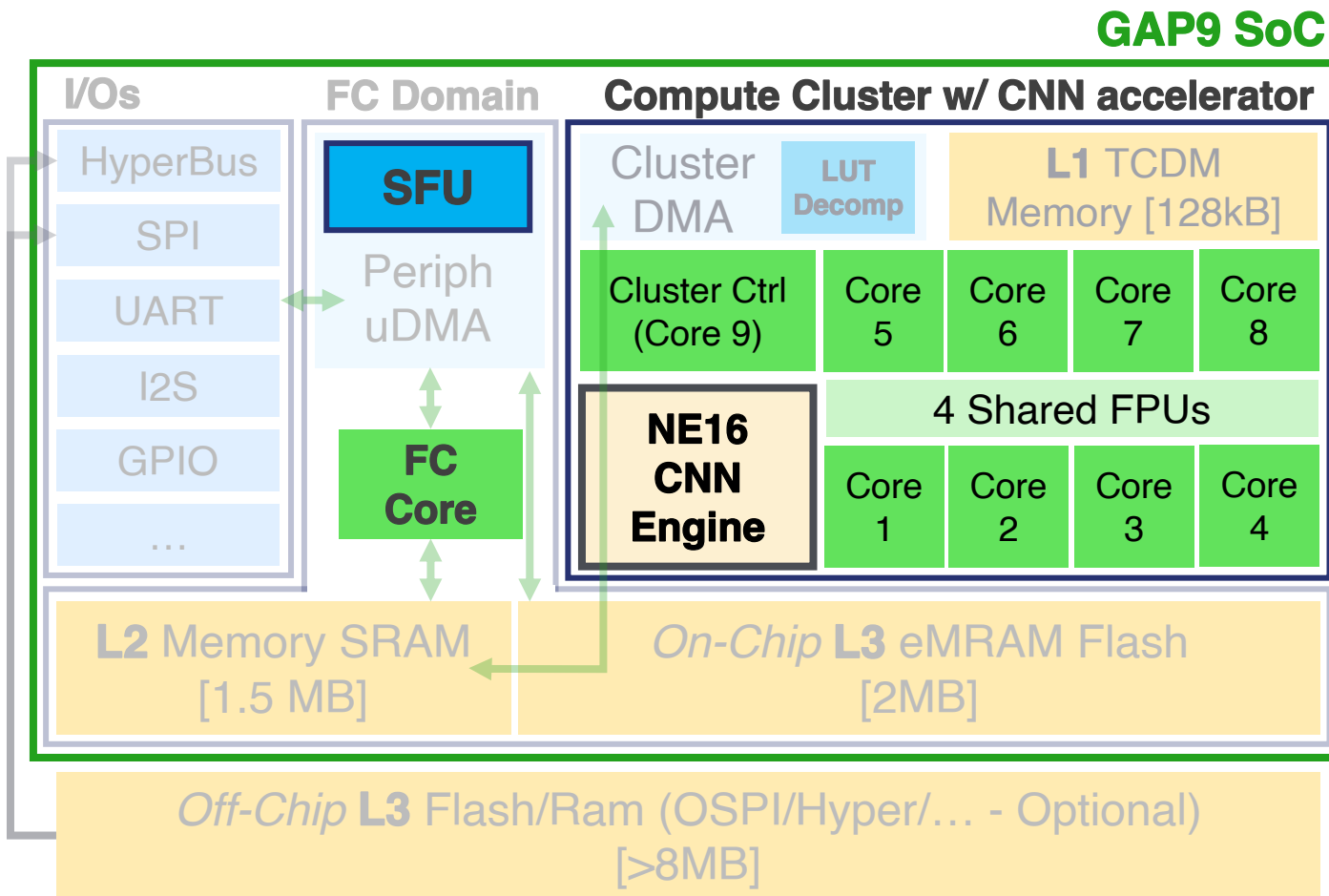
- 4 independent frequency domains: **FC - I/Os - Cluster - SFU**
- “Turn-on when you need”

### Hierarchical Memory Architecture (w/o D-Cache)

- L1: 128kB – 1 Cyc/Access
- L2: 1.5MB – 10-100 Cyc/Access
- L3: >2MB – 100-1000 Cyc/Access
- DMA and uDMA for background copies+decompression



# GAP9 Architecture



## Hierarchical Compute paradigm

- 4 independent frequency domains: **FC - I/Os - Cluster - SFU**
- “Turn-on when you need”

## Hierarchical Memory Architecture (w/o D-Cache)

- L1: 128kB – 1 Cyc/Access
- L2: 1.5MB – 10-100 Cyc/Access
- L3: >2MB – 100-1000 Cyc/Access
- DMA and uDMA for background copies+decompression

## Heterogeneous Compute Units

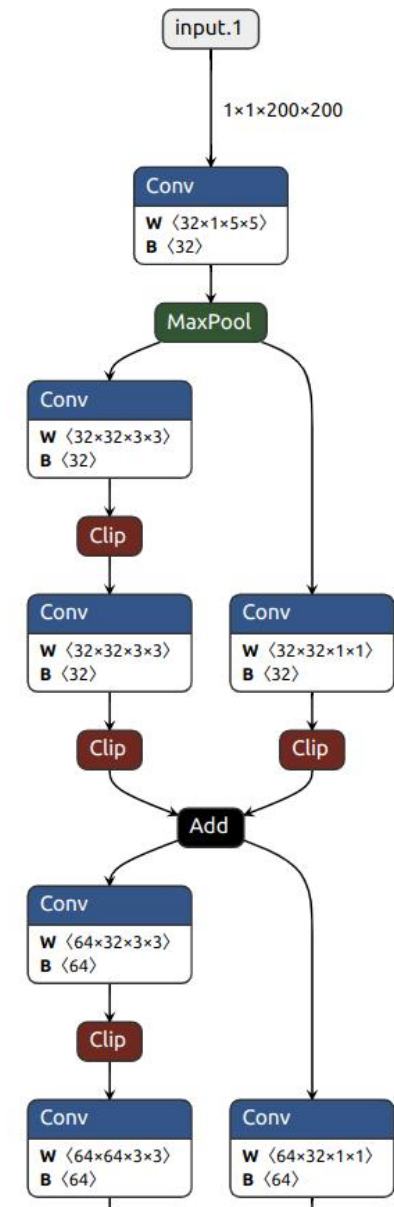
- 10 General Purpose RISC-V Cores
- 4 Shared FPUs (half/single precision)
- Conv/MatMul HW Accelerator (NE16)
- Low-Latency time-domain DSP Accelerator (SFU)

# GAPflow

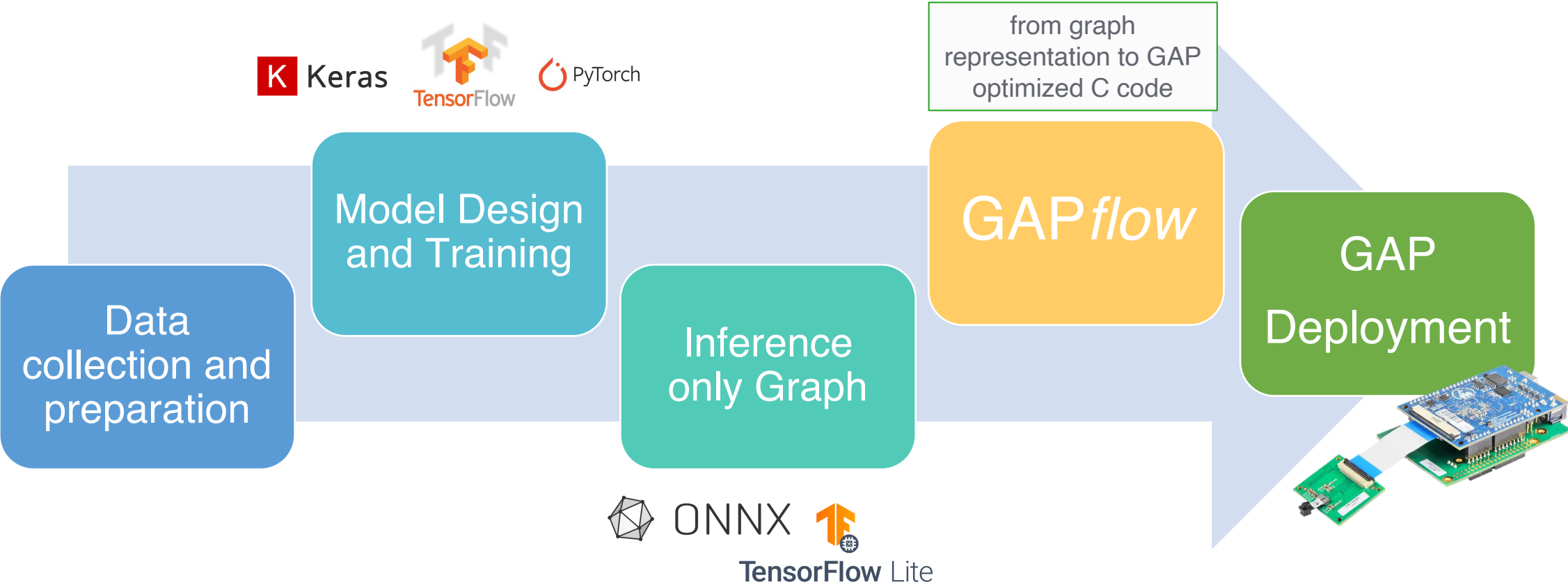
# Neural Network from GAPflow POV

## Computational Graph:

- **Nodes = Computational Layers**
  - e.g. Conv2D, MatMul, MatAdd, Pooling, ...
  - **COMPUTATIONAL COST**
- **Edges = Tensors**
  - e.g. Input, Output, Weights, ...
  - **MEMORY COST**
    - *Static*: constant at every network run (weights/bias)
    - *Dynamic*: different depending on the network inputs (in/out)

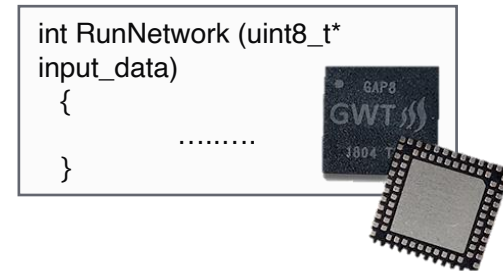
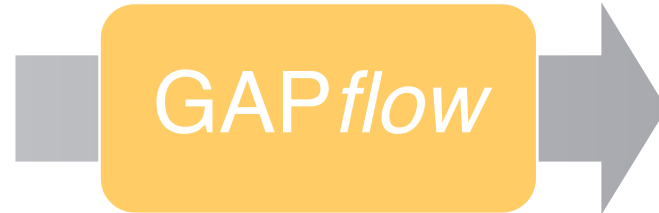
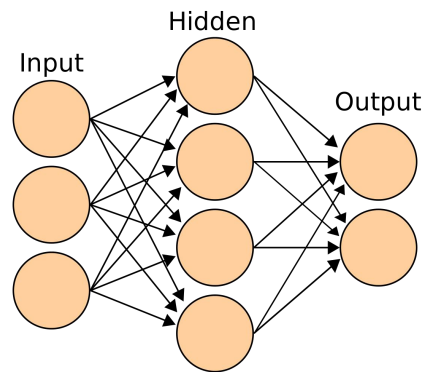


# GAPflow enables DNN inference on Parallel-Ultra Low Power GAP MCUs



# What is GAPflow?

**MAIN GOAL: turn complex DSP/NN computational graphs into optimized C code for GAP9**



# What is GAPflow?

**MAIN GOAL: turn complex DSP/NN computational graphs into optimized C code for GAP9**

## GAP9 Graph Optimization

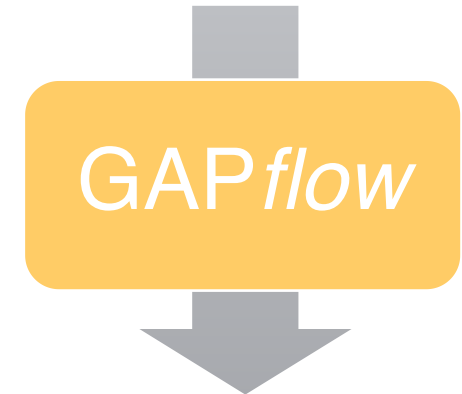
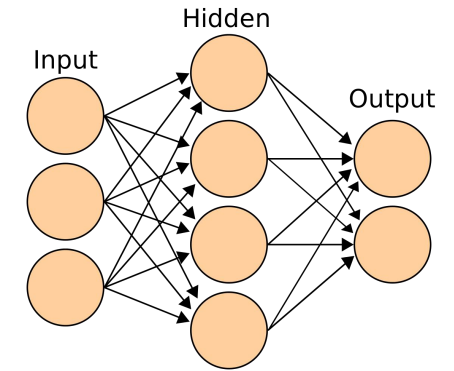
- **Static topology optimization** to minimize number of operations and memory overhead
- **Quantization:** reduce memory usage up to 16x and enable integer only arithmetic when performance is critical

## Validate the Solution

- Validate the numeric precision/accuracy of the deployable model in a user-friendly environment (python)

## Automatic C Code Generation

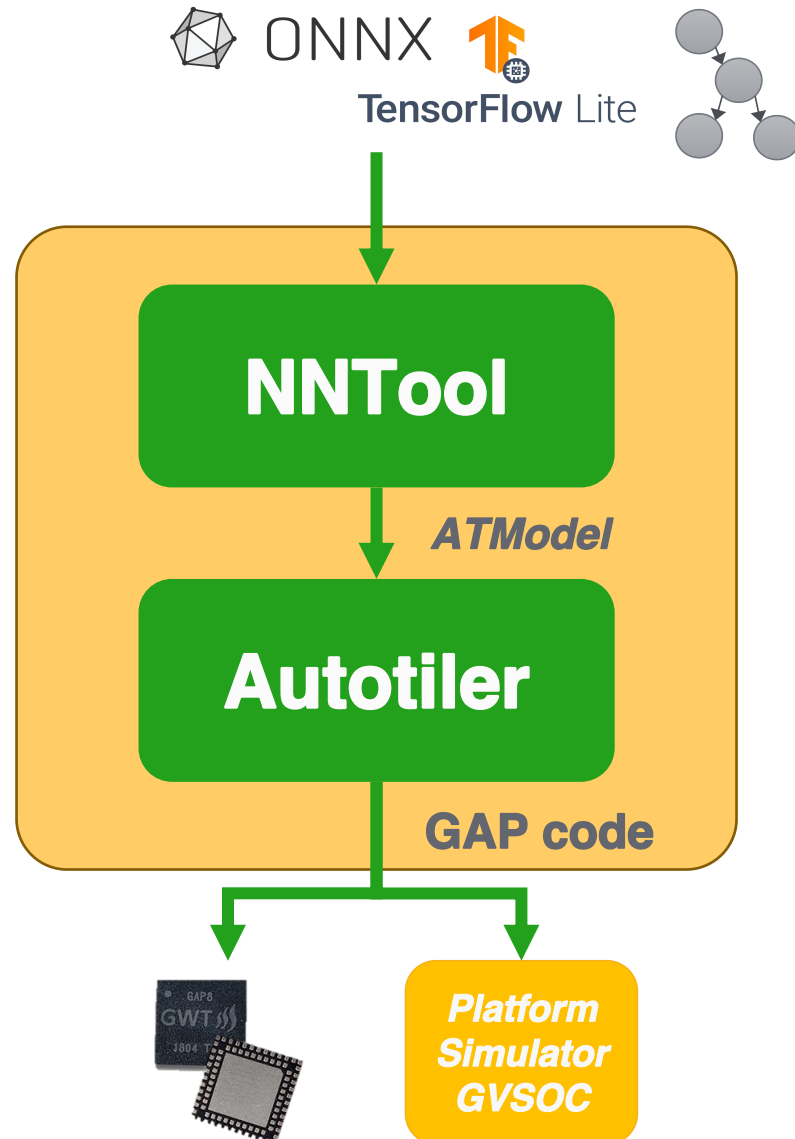
- Map the graph operations into the **Optimized SW library** of GAP9
- **Optimal Memory Management:** automating memory allocation and data transfers



```
int RunNetwork (uint8_t*  
input_data)  
{  
    .....  
}
```

Two images of chips are shown to the right of the code block. The top one is a GAP8 chip with "GAP8" and "GWT" printed on it. The bottom one is a GAP9 chip with "1804" printed on it.

# GAPflow: Overview



## NNTool

- **Static Topology optimizations** (node fusion)
- **Quantization** w/ calibration dataset (optional)
- **Validate** numerically the deployable solution
- Generates an IR of the graph (**ATModel**)

## Autotiler

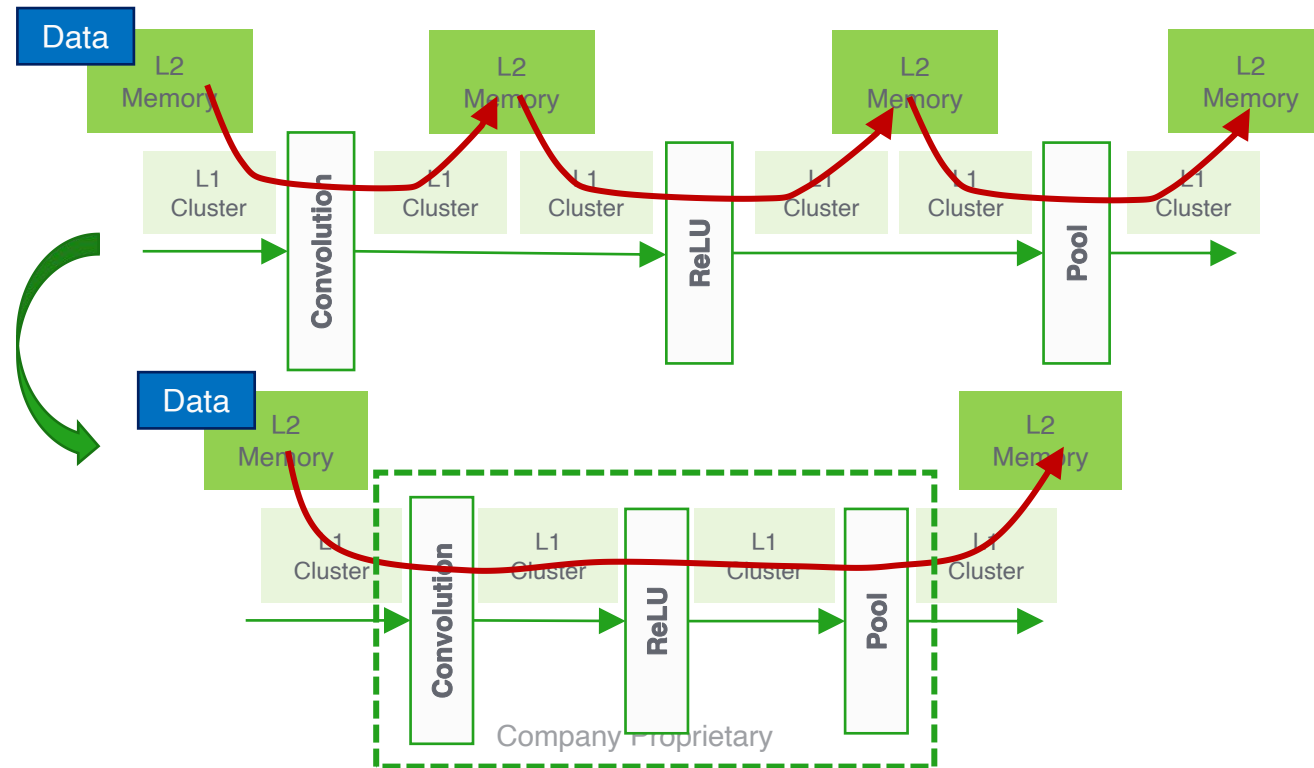
- **Optimizes data movement** across the memory hierarchy
- Computes **optimal tiling sizes**
- Generates **GAP code** with double/triple-buffer mechanism using **optimized SW Library primitives**



# Topology Optimizations

## Minimize number of nodes/edges:

- **Remove** useless reshapes/transpose by moving them accross the graph
- **Layer Fusion:** known sequence of nodes merged together thanks to specialized hand-written backend SW
- **Expressions compiler:** fuses an arbitrary sequence of piecewise/broadcastable operations and dynamically generates GAP C code for it



# Quantization: Background

Convolution Operation

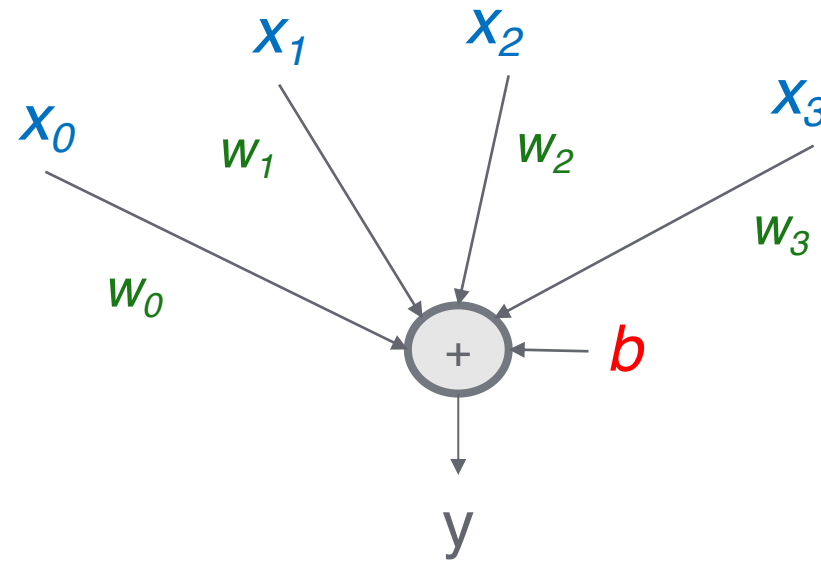
$$\mathbf{y} = \sum \mathbf{x} \cdot \mathbf{w} + \mathbf{b}$$

**x** activation input tensor

**w** weight parameter tensor

**b** bias parameter tensor

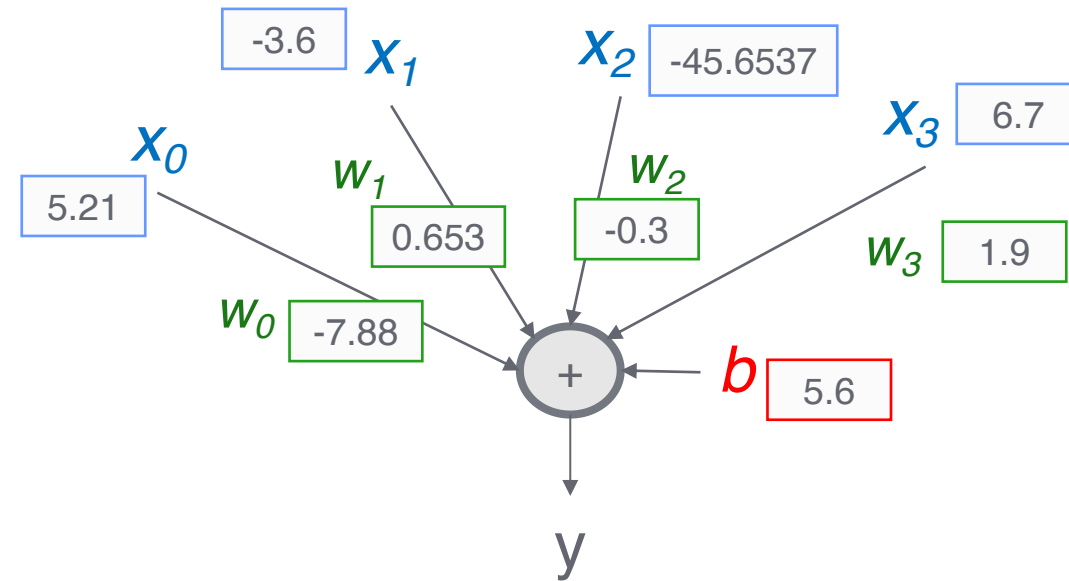
**y** activation output tensor



# Quantization: Background

Convolution Operation

$$y = \sum x \cdot w + b$$



DL frameworks operates with **real numbers**

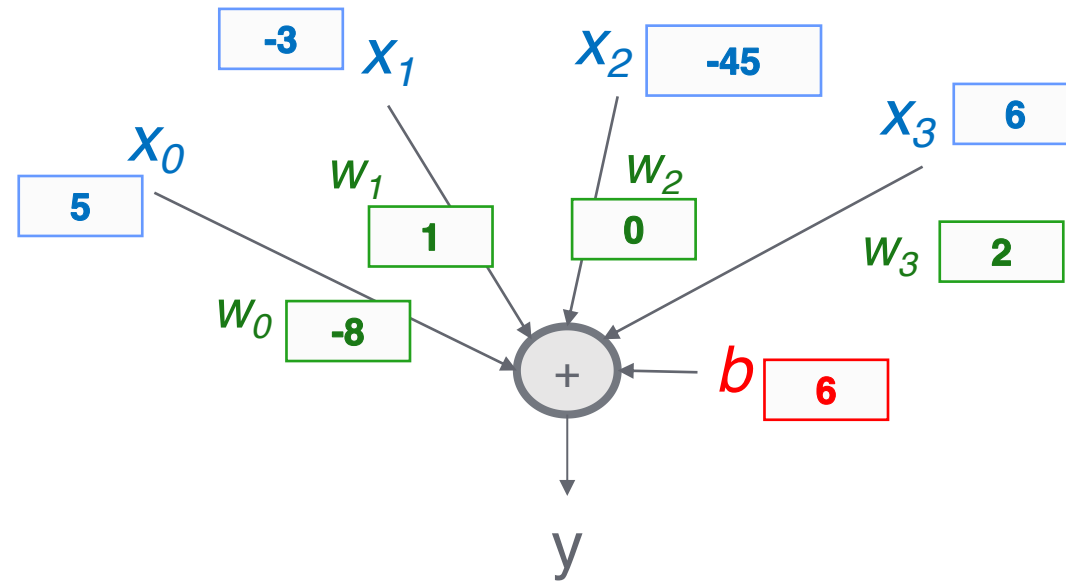
👉 Floating-point 32-bit format (FP32)

👉 Inference requires FPU engines

# Quantization: Background

Convolution Operation

$$y = \sum x \cdot w + b$$



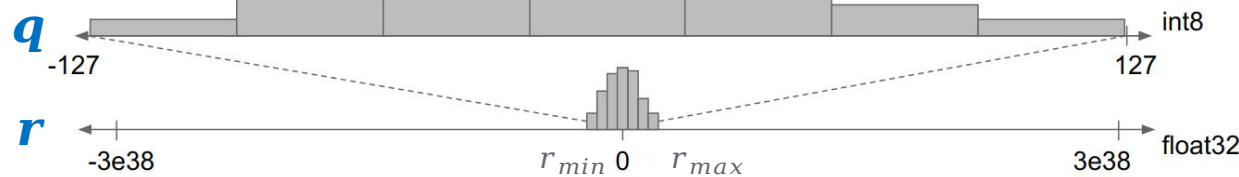
**Quantization** maps any real value into a set of **integer values**

- 👍 **COMPRESSION** to n-bit integer values up to **4x** compression with 8-bit quantization
- 👍 **LATENCY** Inference requests integer low-precision operation 8-bit convolution up to **4x** (even more due to the lower BW) faster than FP32 in SW (even more with dedicated HW like NE16)

# Quantization: Background

- 4x memory reduction
- >4x speed up thanks to specialized HW
- But accuracy???

through an affine



$$r = S (q - Z)$$

$Z = 0$  if symmetric ranges  $r_{max} = -r_{min}$

$$\frac{r_{max} - r_{min}}{2^{nbits}}$$

## Convolution Operation

$X$  is a quantized value,  $x$  is a real value

$$y = \sum x \cdot w \longrightarrow S_y Y = \sum S_x X \cdot S_w W$$

Affine transformation

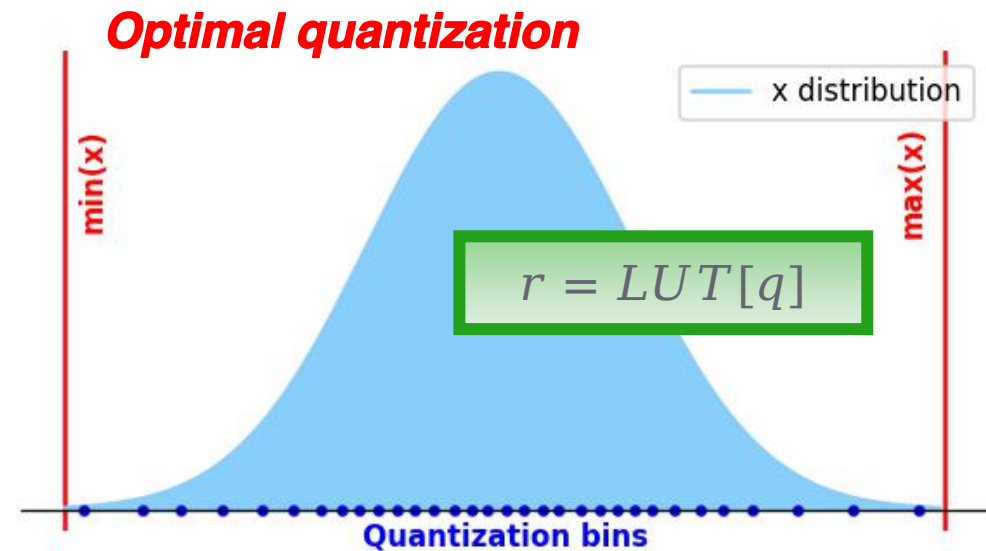
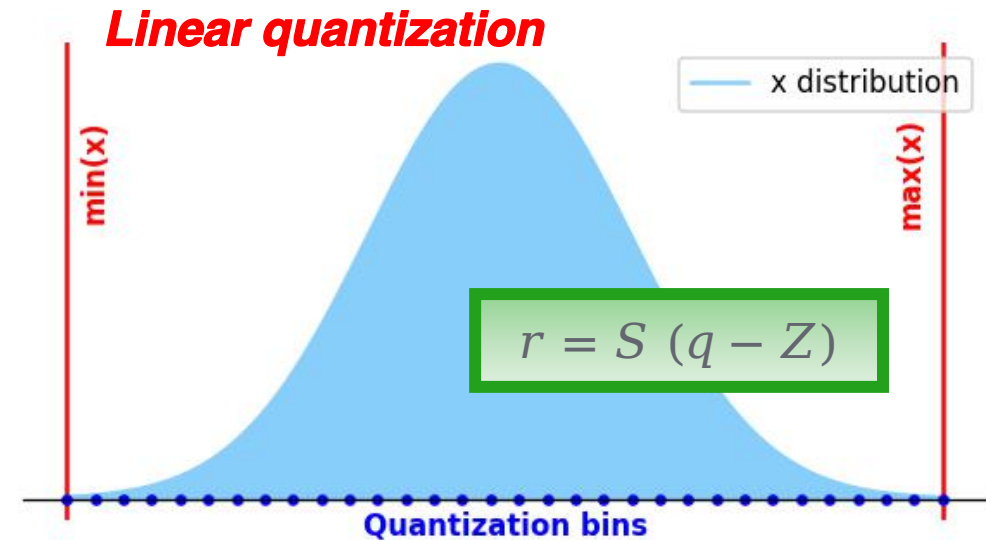
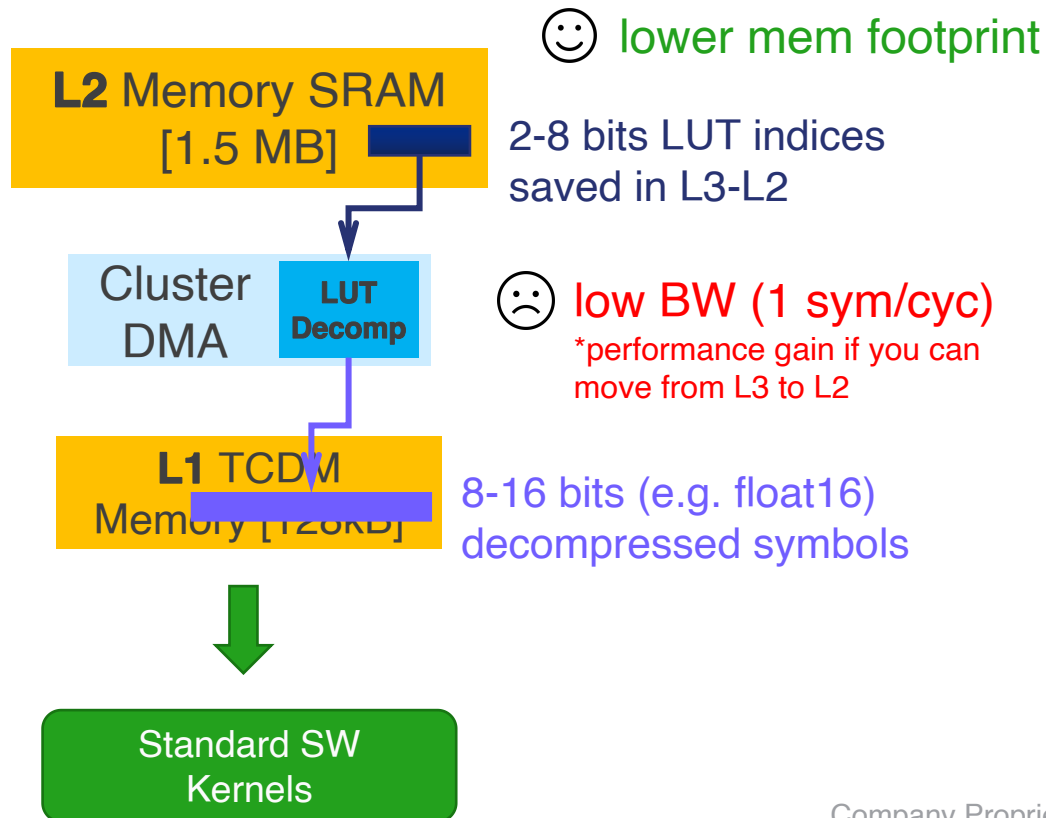
$$Y = \frac{S_x S_w}{S_y} \sum X \cdot W$$

Integer only MACs

INT8

# Optimal quantization: Decompressor

- LUT based quantization allows non-uniform quantization
- Clustering (KMeans) can lead to better approximation (higher compression rates)



# Quantization

$$r = S (q - Z)$$

$\frac{r_{max} - r_{min}}{2^{nbits}}$        $Z = 0$  if symmetric ranges  $r_{max} = -r_{min}$

Scheme	Weights	Activations	Affine Transformation
<b>POW2 (QX.Y)</b>	16-bit (symmetric, PerTensor)	16-bit (symmetric, PerTensor)	$r = q * 2^{-N}$ (1)
<b>int8</b>	2-8bits (symmetric, PerChannel)	8- or 16bit (asym. PerTensor)	$r = (q - Z) * S$ (1)
<b>Float16</b>	16-bit (bfloat16/f16ieee)	16-bit (bfloat16/f16ieee)	- (2)
<b>LUT</b>	2-8bits (LUT)	(not supported)	$r = LUT[q]$

(1) Quantization Ranges obtained:

- From a graph quantized with third-parties tools: *Onnx+NNCF* or *TFLite* quantization
- Using GAP Nntool Post-Training Quantization with a calibration dataset

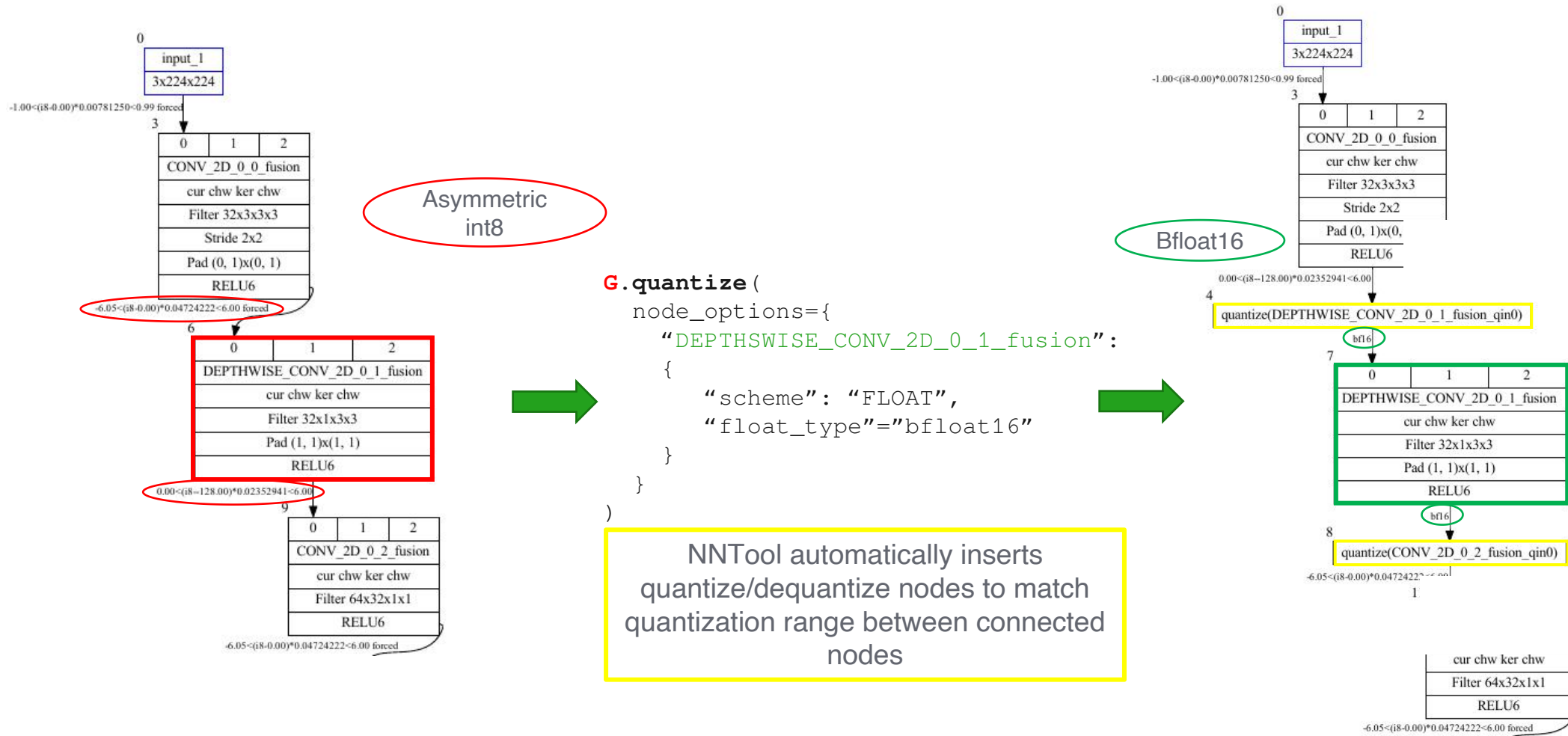
(2) Cast from FP32 to FP16 done by GAP NNtool. No Calibration dataset required. Typically, lossless.

\* TFLite quantization, <https://arxiv.org/abs/1712.05877>

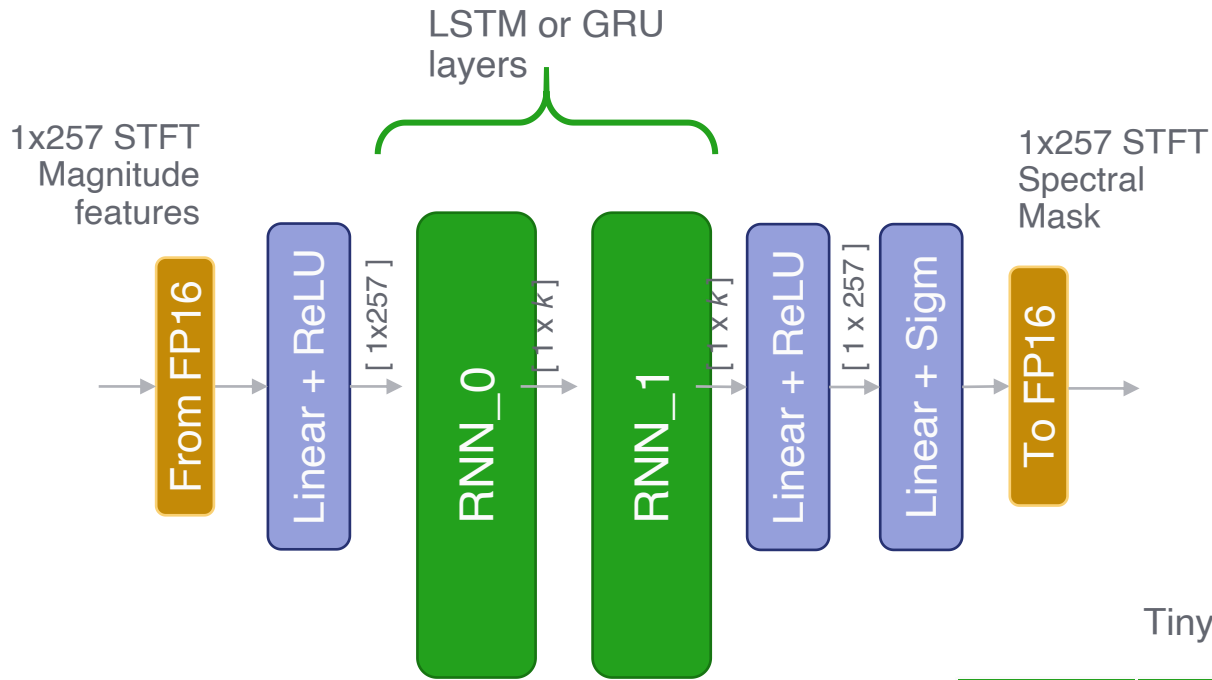


# Mixed-Precision Quantization

GAP NNTool enables layer-wise selection of quantization scheme and number of bits



# Mixed-Precision Quantization: Use case



	LSTM256	GRU256	LSTM128	GRU128
<b>k</b>	256	256	128	128
<b>RNN_0</b>	LSTM(257,256)	GRU(257, 256)	LSTM(257,128)	GRU(257, 128)
<b>RNN_1</b>	LSTM(257,256)	GRU(257, 256)	LSTM(128, 128)	GRU(128, 128)
<b>Param</b>	1.24 M	0.985 M	0.493 M	0.411 M
<b>% rnn params</b>	84%	80%	66.5%	59.8%

The majority of the weights are due to RNN layers!

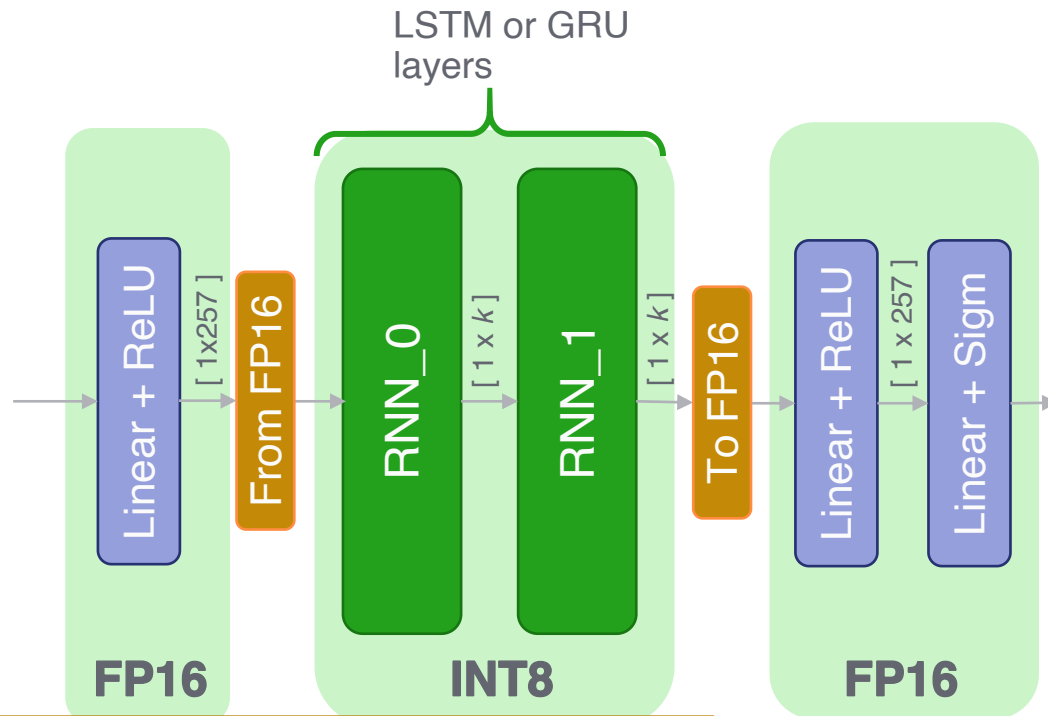
TinyDenoiser models trained on Valentini dataset: FP32 baseline

Quant Type	LSTM256			GRU256			LSTM128			GRU128		
	PESQ	STOI	Mem	PESQ	STOI	Mem	PESQ	STOI	Mem	PESQ	STOI	Mem
<b>FP32</b>	<b>2.79</b>	<b>0.94</b>	4.75	<b>2.78</b>	<b>0.94</b>	3.76	<b>2.76</b>	<b>0.94</b>	1.88	<b>2.69</b>	<b>0.94</b>	1.56
<b>FP16</b>	<b>2.79</b>	<b>0.94</b>	2.37	<b>2.78</b>	<b>0.94</b>	1.88	<b>2.76</b>	<b>0.94</b>	0.94	<b>2.69</b>	<b>0.94</b>	0.78
<b>INT8</b>	<b>2.42</b>	<b>0.92</b>	1.18	<b>2.48</b>	<b>0.93</b>	0.93	<b>2.51</b>	<b>0.92</b>	0.47	<b>2.36</b>	<b>0.93</b>	0.39

**INT8** is lightweight but **LOSSY**

- avg PESQ loss: **-0.3**, STOI loss: **0.015**
- **2x** mem compression

# Mixed-Precision Quantization: Use case



**No need for expensive QAT !!!**

**MixedFP16-INT8** present low accuracy degradation (**PESQ: 0.06, STOI: <0.01**) while **1.4-1.7x** mem compression vs FP16

	LSTM256	GRU256	LSTM128	GRU128
<b>k</b>	256	256	128	128
<b>RNN_0</b>	LSTM(257,256)	GRU(257, 256)	LSTM(257,128)	GRU(257, 128)
<b>RNN_1</b>	LSTM(257,256)	GRU(257, 256)	LSTM(128, 128)	GRU(128, 128)
<b>Param</b>	1.24 M	0.985 M	0.493 M	0.411 M
<b>% rnn params</b>	84%	80%	66.5%	59.8%

The majority of the weights are due to RNN layers!

TinyDenoiser models trained on Valentini dataset: FP32 baseline

Quant Type	LSTM256			GRU256			LSTM128			GRU128		
	PESQ	STOI	Mem	PESQ	STOI	Mem	PESQ	STOI	Mem	PESQ	STOI	Mem
<b>FP32</b>	<b>2.78</b>	<b>0.94</b>	4.75	<b>2.78</b>	<b>0.94</b>	3.76	<b>2.76</b>	<b>0.94</b>	1.88	<b>2.69</b>	<b>0.94</b>	1.56
<b>FP16</b>	<b>2.78</b>	<b>0.94</b>	2.37	<b>2.78</b>	<b>0.94</b>	1.88	<b>2.76</b>	<b>0.94</b>	0.94	<b>2.69</b>	<b>0.94</b>	0.78
<b>INT8</b>	<b>2.42</b>	<b>0.92</b>	1.18	<b>2.48</b>	<b>0.93</b>	0.93	<b>2.51</b>	<b>0.92</b>	0.47	<b>2.36</b>	<b>0.93</b>	0.39
<b>MixFP16-INT8</b>	<b>2.73</b>	<b>0.93</b>	1.37	<b>2.72</b>	<b>0.94</b>	1.13	<b>2.69</b>	<b>0.93</b>	0.67	<b>2.63</b>	<b>0.94</b>	0.55

# Fast prototyping

- Deploy in few lines
- Complete control of:
  - Graph manipulations
  - Quantization
  - Testing
  - Deployment settings
- Check of consistency

## *Prepare the graph for deployment*

```
G = NNGraph.load_graph("model.onnx")
G.adjust_order()
G.fusions("scaled_match_group")
stats = G.collect_statistics(repr_dataset())
G.quantize(stats, graph_qopts, node_qopts)
```

## *Test the final graph in python*

```
for data in test_data:
    outs = G.execute(data, dequantize=True)
    ok = check(outs)
```

## *If not satisfied: requantize*

```
if not ok:
    # update node_qopts and start over
```

## *Check consistency on target*

```
res = G.execute_on_target(data, model_settings)
check_equal(res.output_tensors, outs)
res.print_basic_mem_infos()
res.print_performance()
...
```

# Validation of the solution

- Check deployed model accuracy:

☹️ **On-device**: use directly the final platform to test the accuracy (**very slow**)

😊 **NNTool**: bit-accurate numpy backend, the user can test accuracy in a python environment without need of device (**fast**)

```
G = NNGraph.load_graph(file_path)
stats = G.collect_statistics(calibration_dataset)
G.quantize(stats, quantization_options)
G.adjust_order()
G.fusions("scaled_match_group")
# Ready for inference
acc1 = 0
for in_data, target in test_dataset:
    outq = G.execute(in_data, quantize=True)
    acc1 += np.argmax(outq[-1][0]) == target
```

**Prepare your model for deployment**  
(quantization+graph manipulation)

**Run inference and check results**

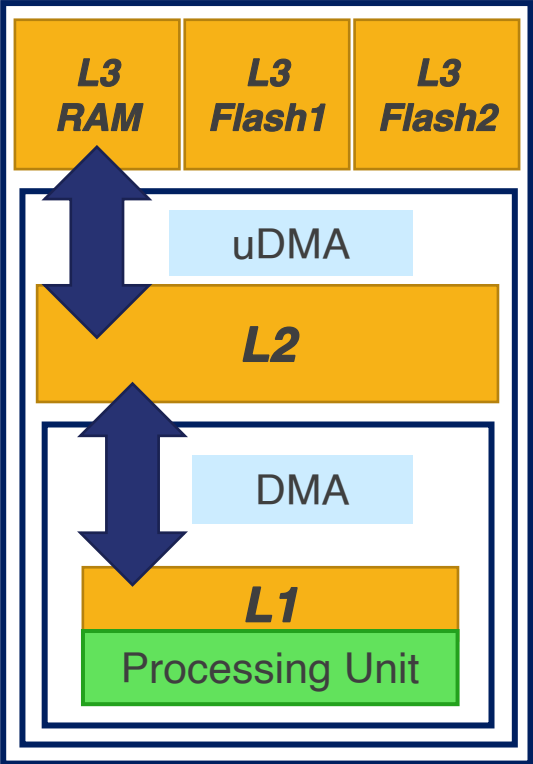
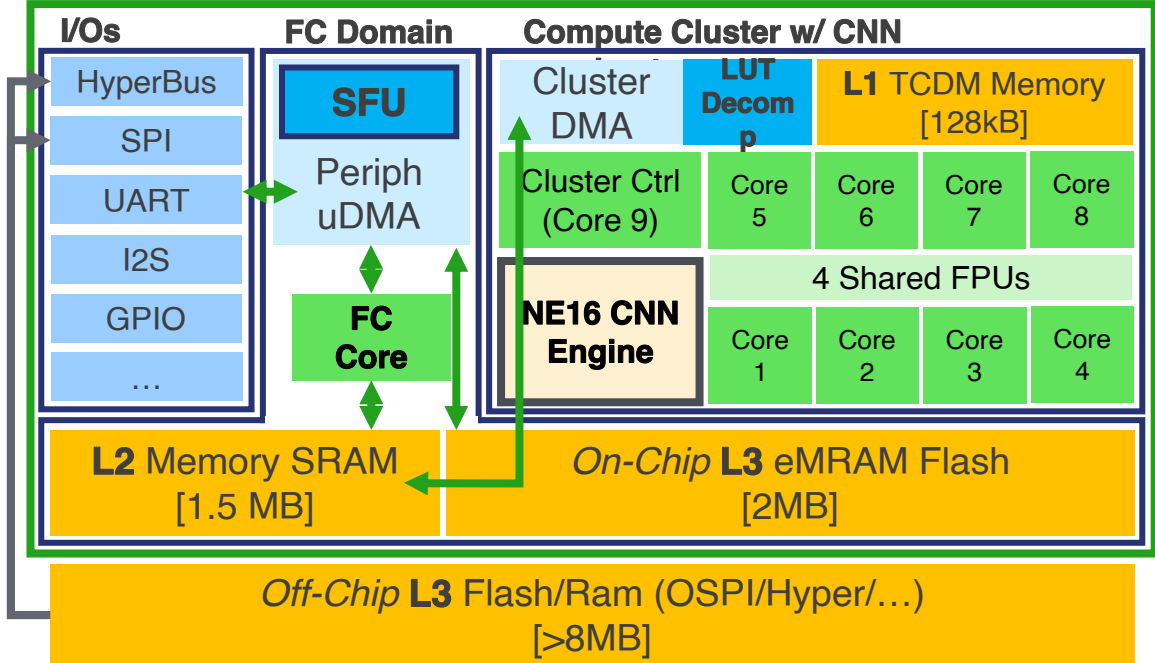
# GAP Autotiler

GAP does not have a data CACHE:

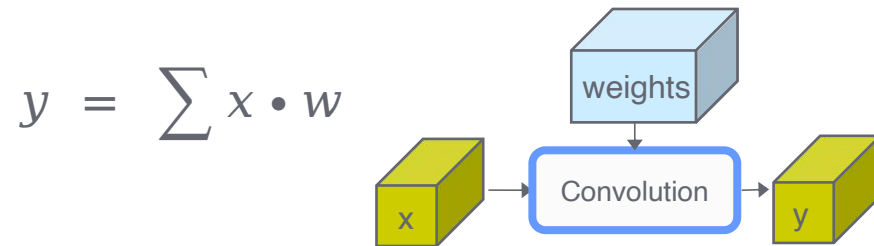
- Silicon Area
- Energy Efficiency
- NN/DSP algo have predictable data traffic



Generate C code for all data movement at compile time



# Autotiler User Kernel: NN Node to the GAP architecture



## Computation dataflow

Ahead of time

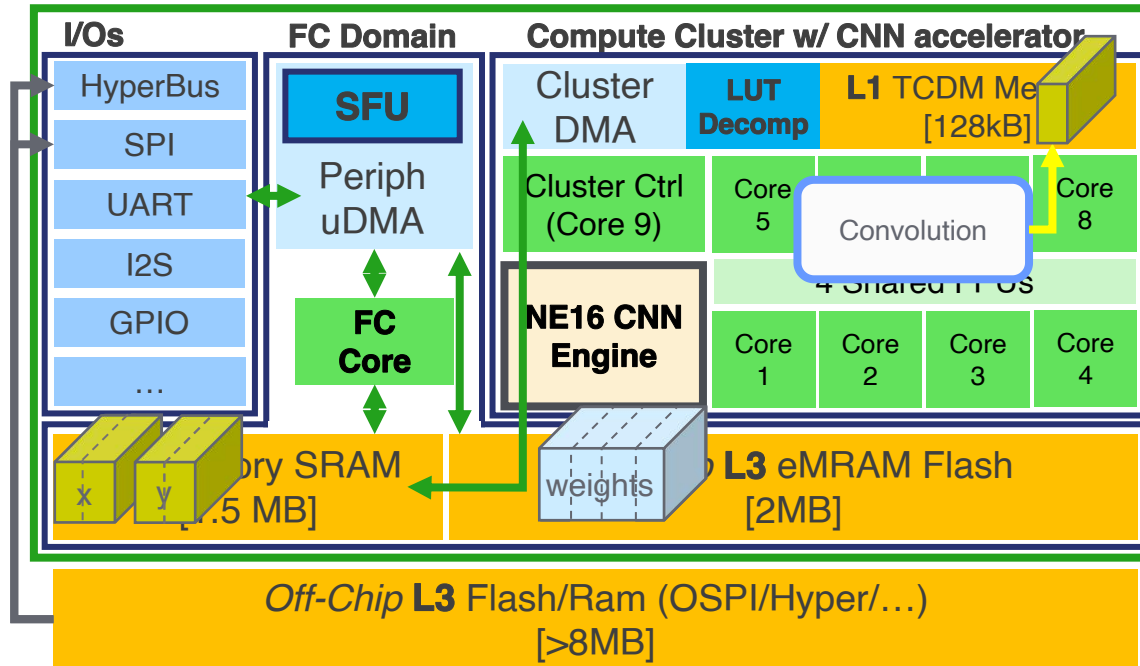
Store data (parameters & input vector) in L2 (or L3)

At run time, for any computational node:

Partition and Load data (parameters & input tensors) to L1

Run data-parallel/CNN Engine computation

Store data (output tensors) back in L2 (or L3)



```
static void Conv_Layer0
(
    signed char * In,      // input L2 vector
    signed char * Weights, // input L2 vector
    signed char * Bias,   // input L2 vector
    signed char * Out,    // output L2 vector
){
    //tile sizes of In, Weights, Bias computed offline
    //L1 buffer allocated to handle double buffering
    // two L1 memory buffers for double buffering

    DMA load first tiles to L1 memory buffer

    for any tile of In, Weights, Bias tensors:

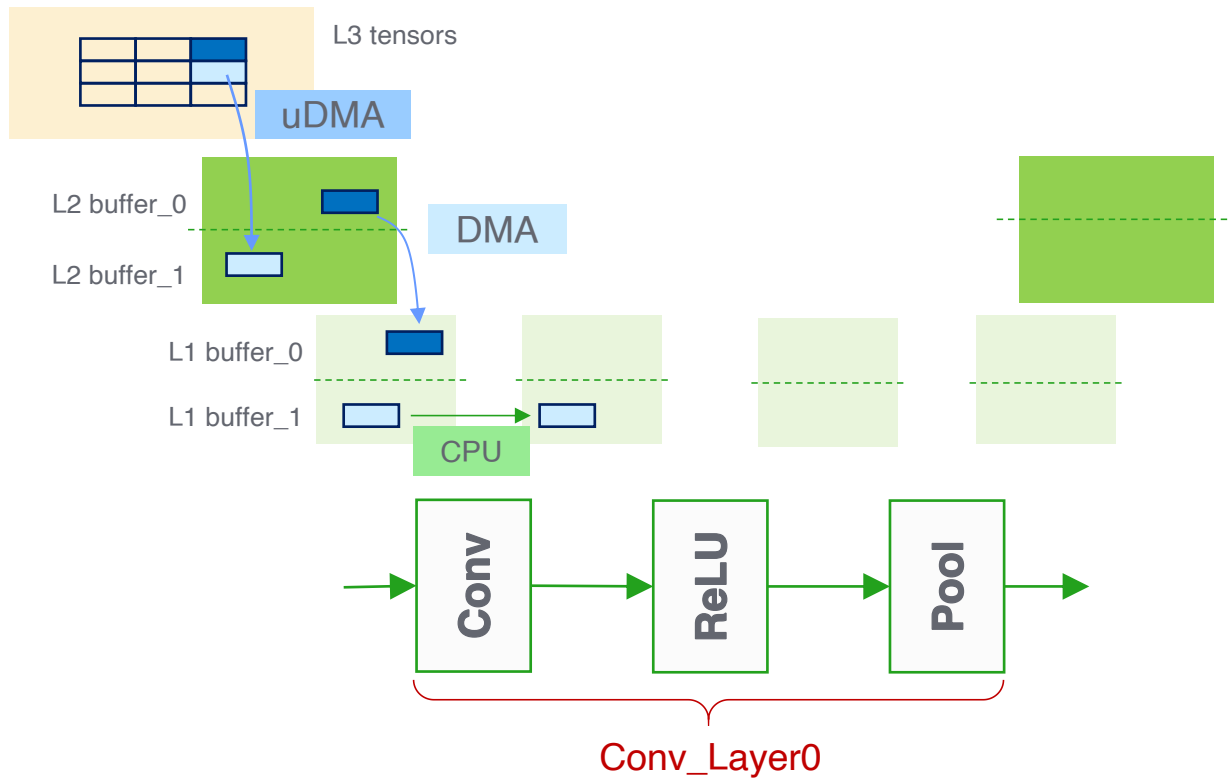
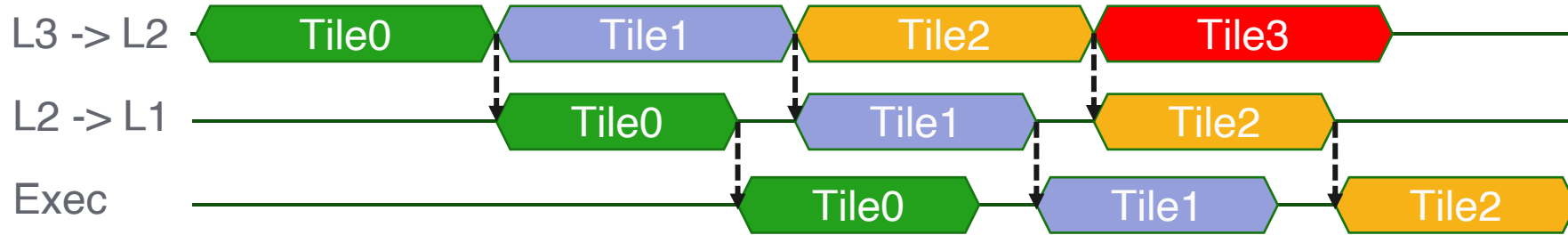
        DMA load next tiles to L1 memory buffer

        ParConv() on L1 tile
        ParReLU() on L1 tile
        ParPool() on L1 tile

        DMA write results (Out) to L2
}
```



# Mapping a NN Node to the GAP HW/SW architecture



```

static void Conv_Layer0
{
    signed char * In,      // input L3 vector
    signed char * Weights, // input L3 vector
    signed char * Bias,   // input L3 vector
    signed char * Out,    // output L3 vector
}{
    //tile sizes of In, Weights, Bias computed offline
    //L1 buffer allocated to handle double buffering
    // two L1 memory buffers for double buffering

    uDMA load first tiles to L2 memory buffer

    DMA load first tiles to L1 memory buffer

    for any tile of In, Weights, Bias tensors:

        uDMA load next next tiles to L2 memory buffer

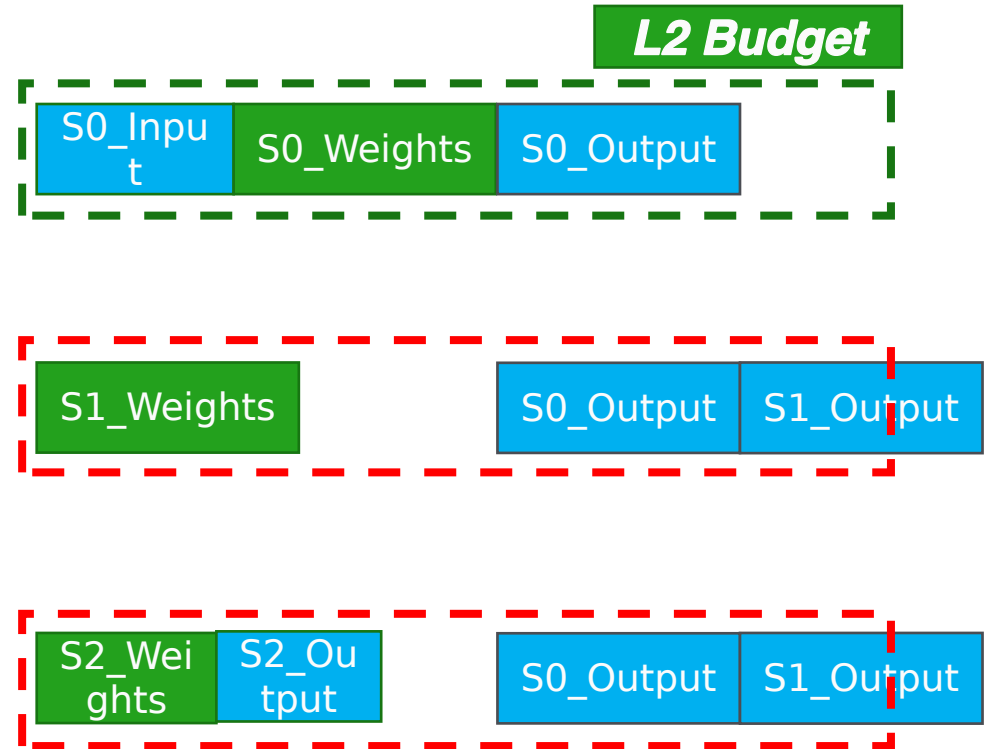
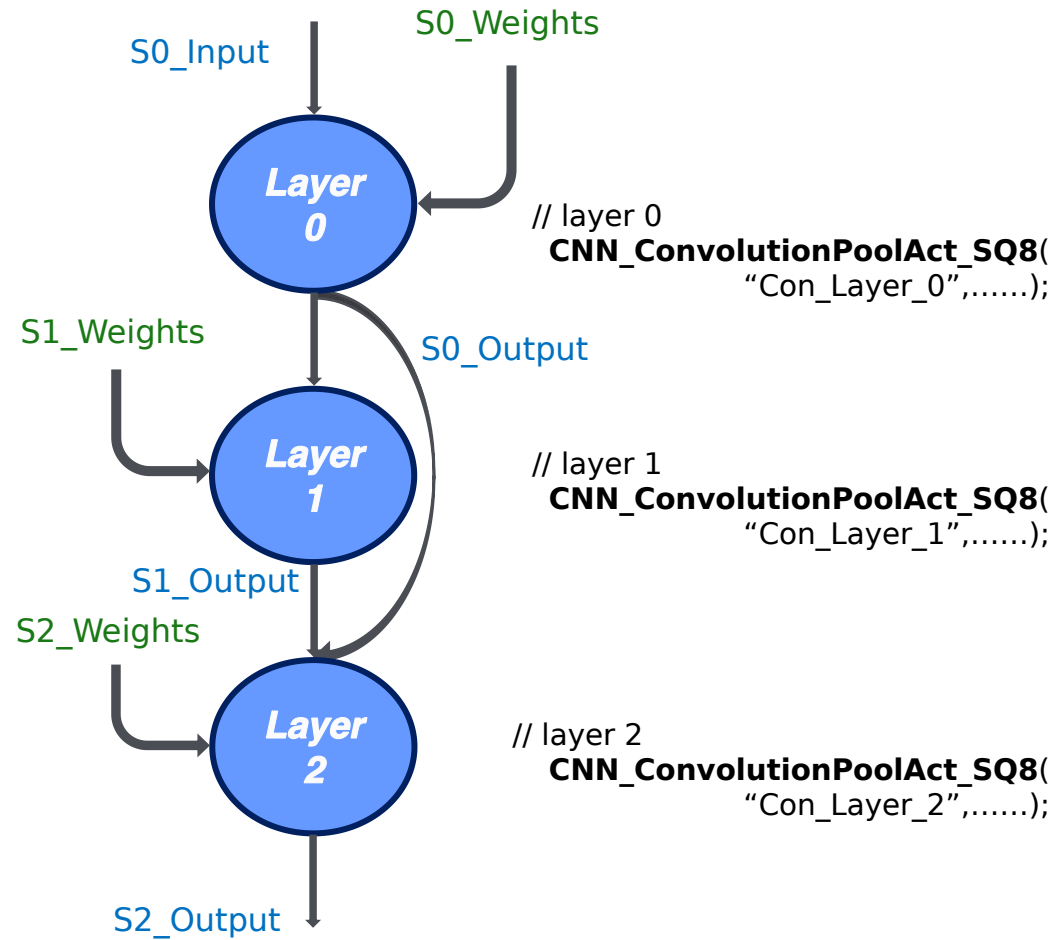
        DMA load next tiles to L1 memory buffer

        ParConv() on L1 tile
        ParReLU() on L1 tile
        ParPool() on L1 tile

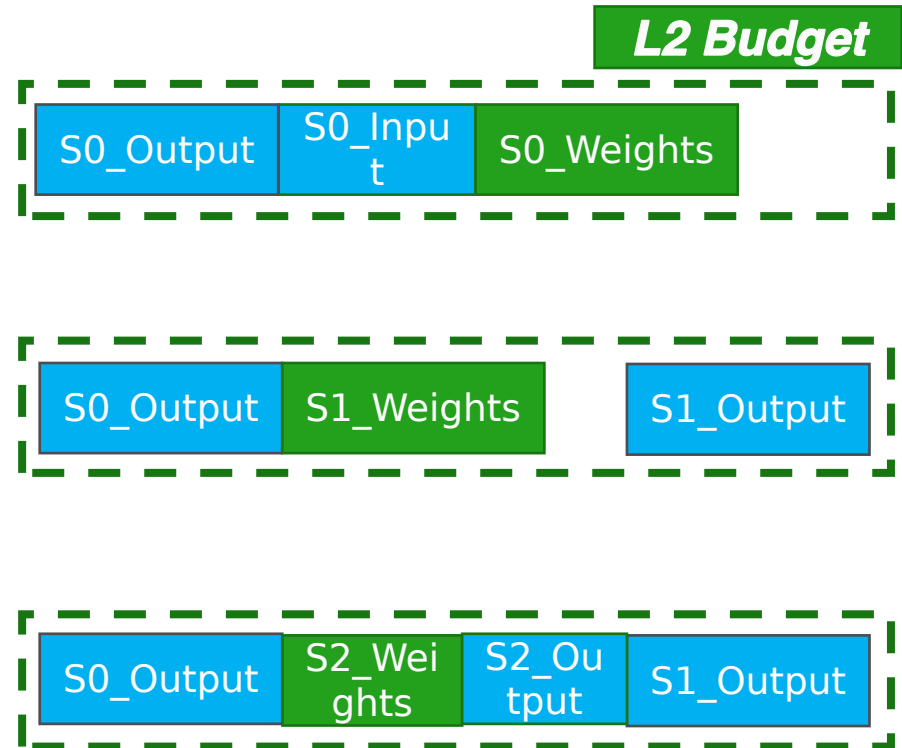
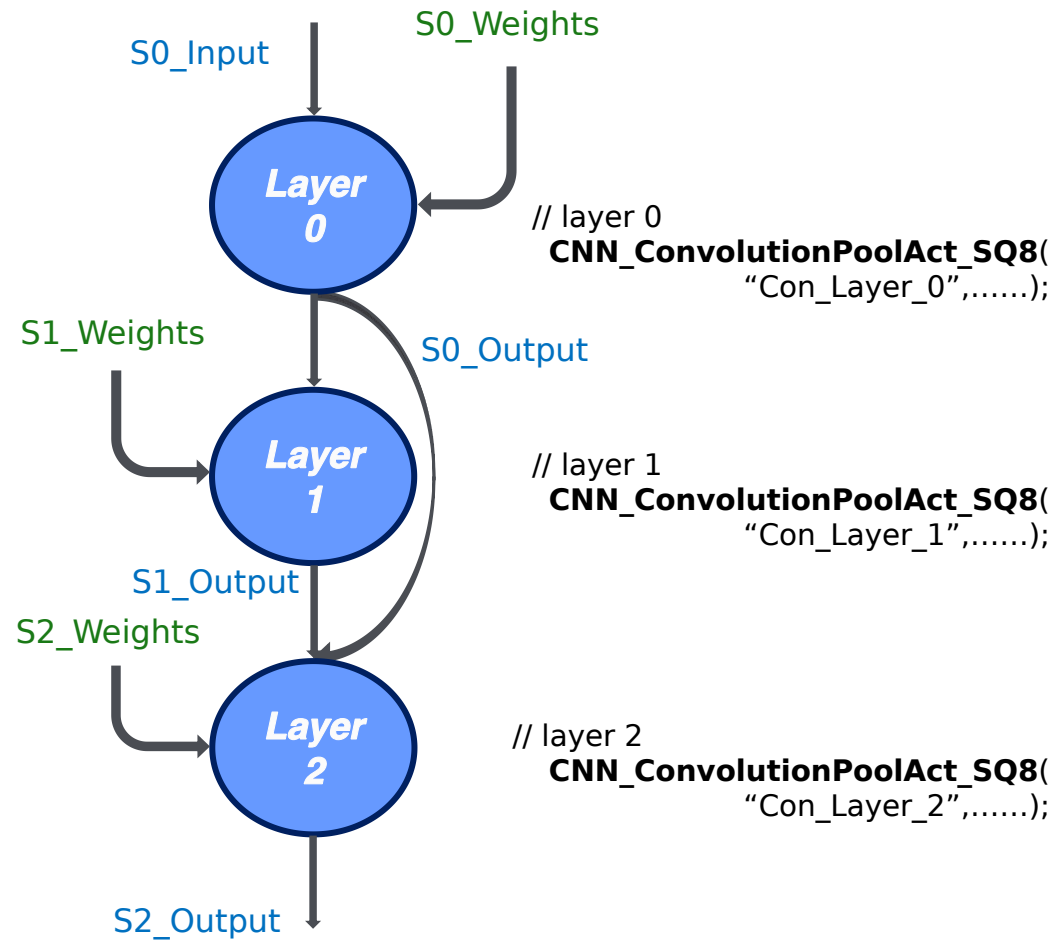
        DMA write results (Out) to L2

        uDMA write prev results to L3
}
    
```

# Autotiler Graph: Static allocation

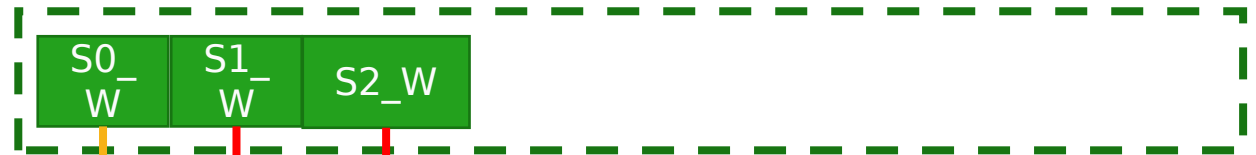
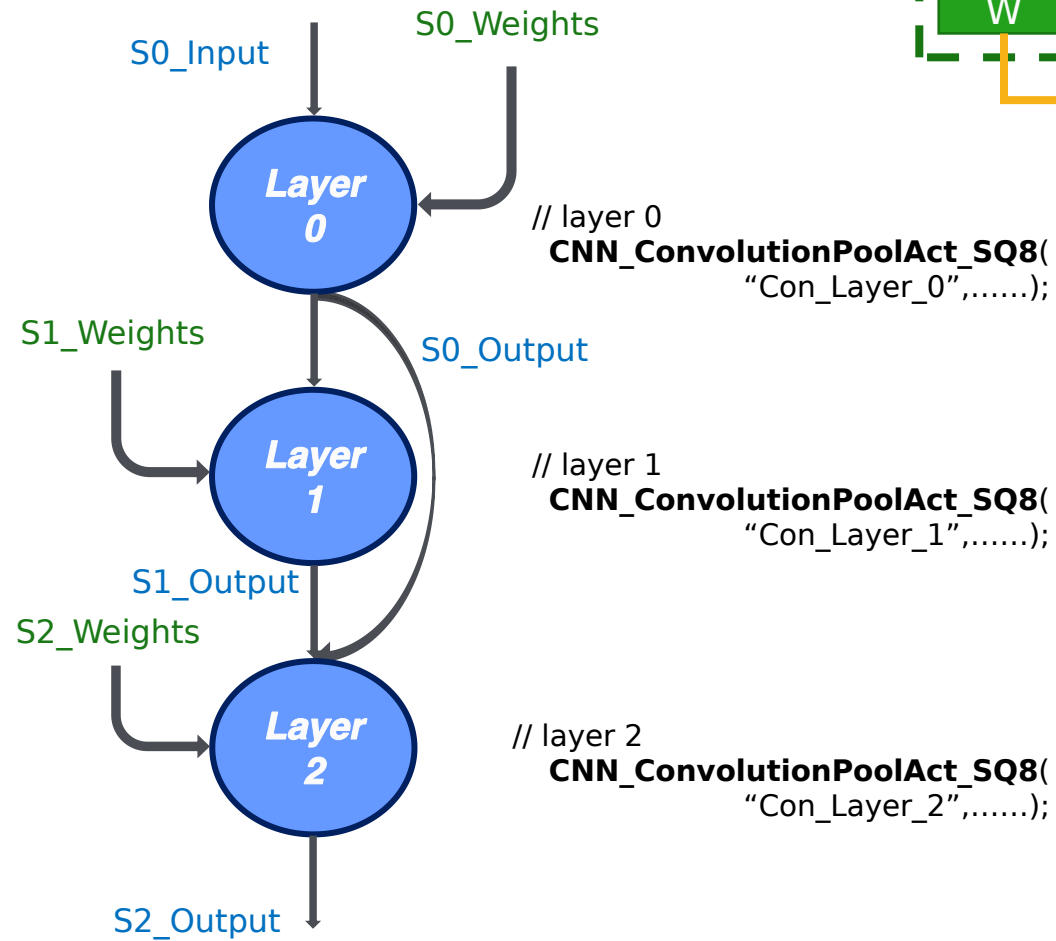


# Autotiler Graph: Static allocation

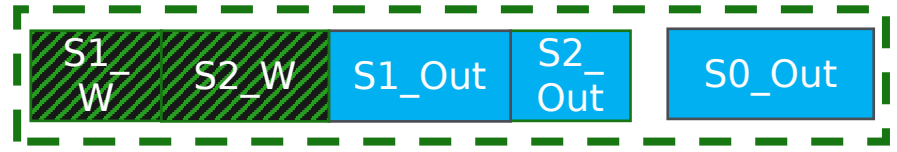


# Autotiler Graph: Static allocation (tensor promotion)

**L3 Budget**

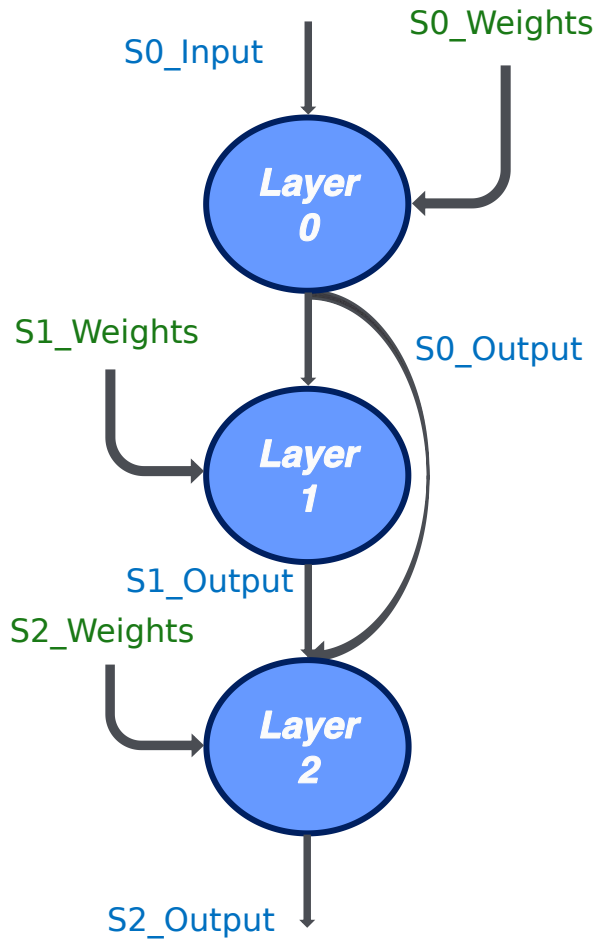


**L2 Budget**

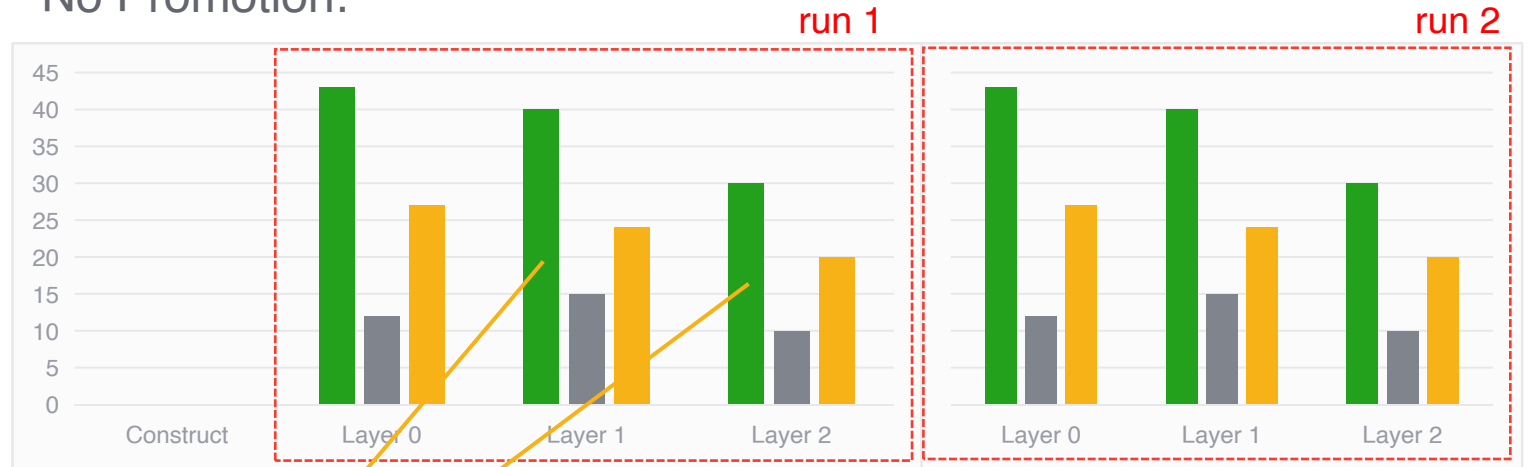


# 2. Graph mode: Mem allocation inter-layer

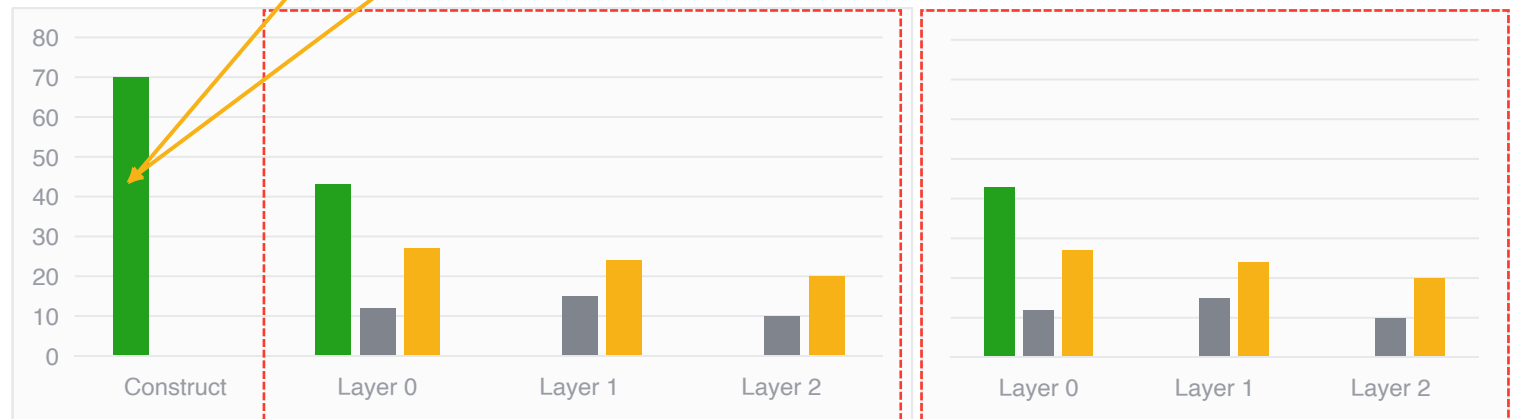
- Energy Weights L3->L2
- Energy L2->L1
- Energy Execution



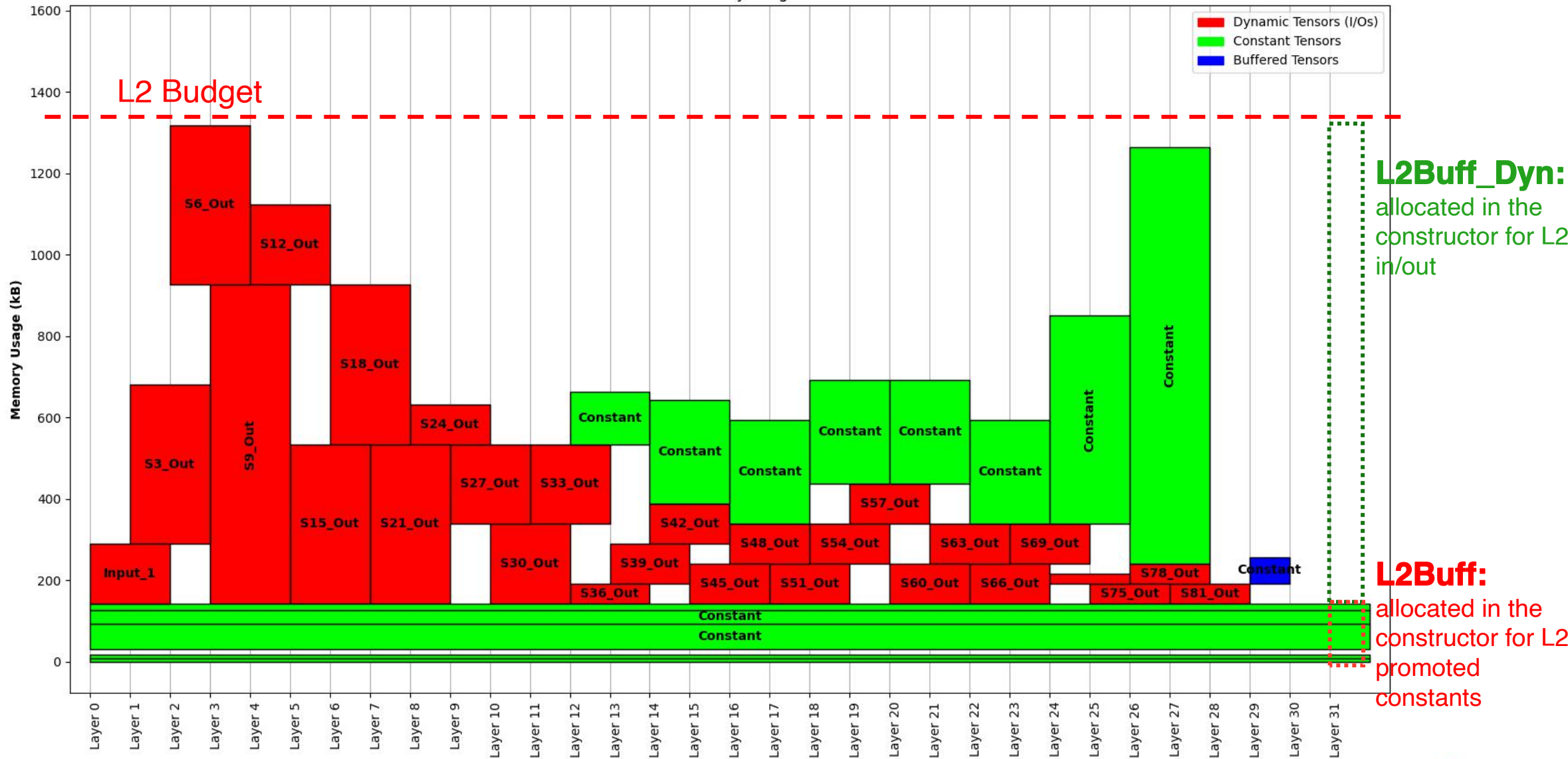
No Promotion:



Static Promotion of constants

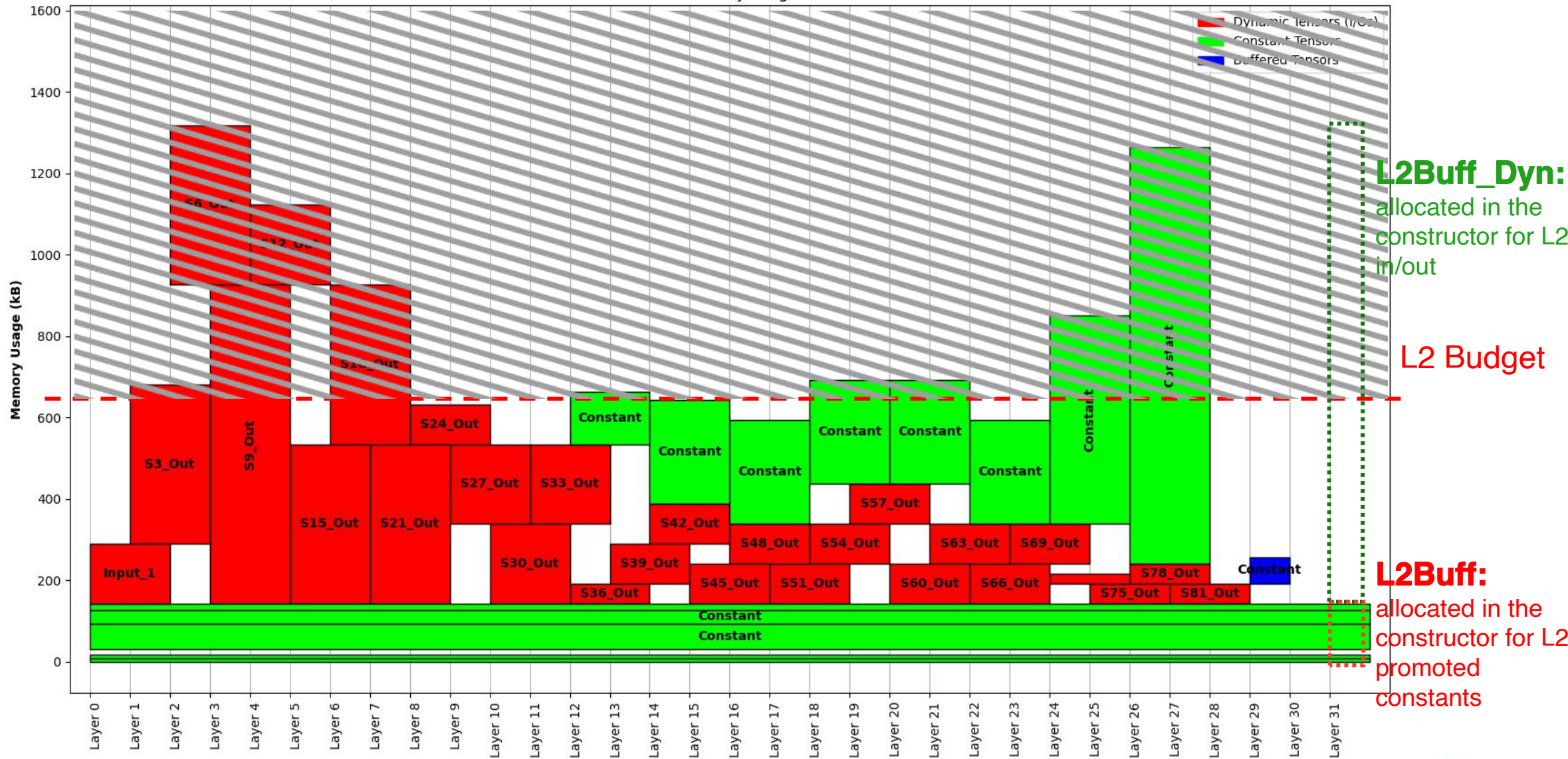


# MobilenetV1 L2 RAM Graph Memory allocation Strategy

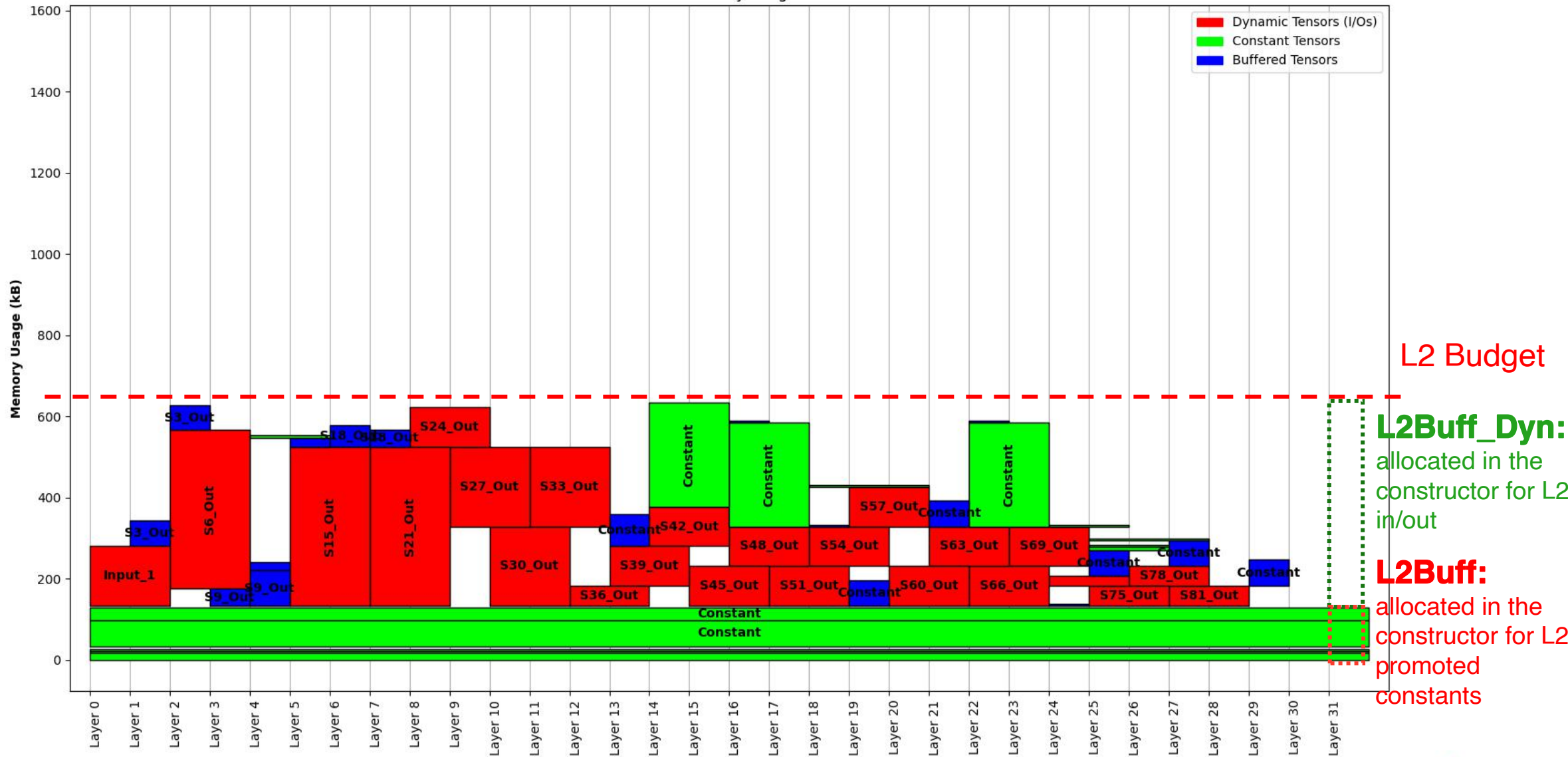




# MobilenetV1 L2 RAM Graph Memory allocation Strategy



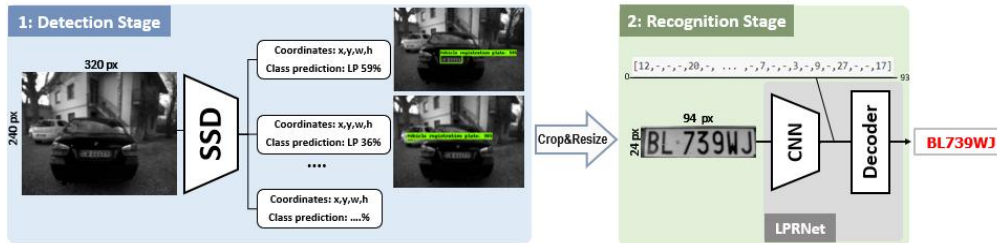
# MobilenetV1 L2 RAM Graph Memory allocation Strategy





# Enabling complex NN on energy constrained devices

## Multi NN – Licence Plate



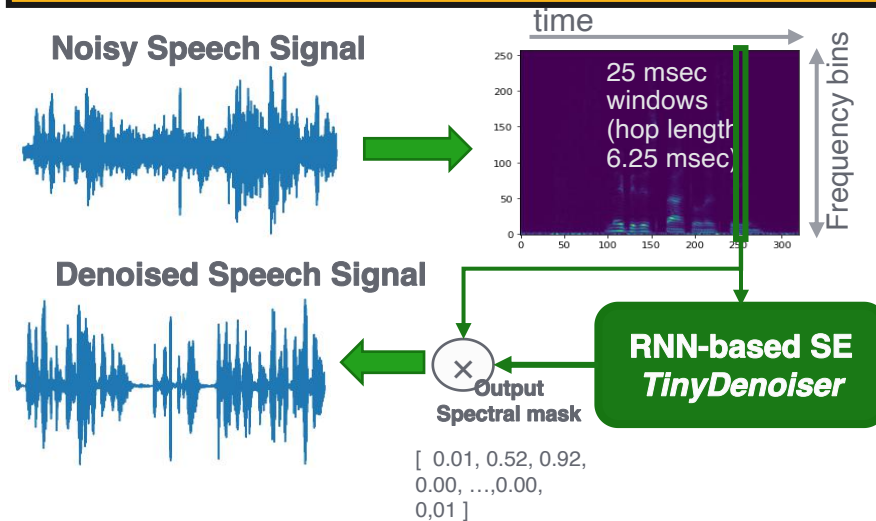
<https://ieeexplore.ieee.org/abstract/document/9401730>

## Smart glasses - object detection



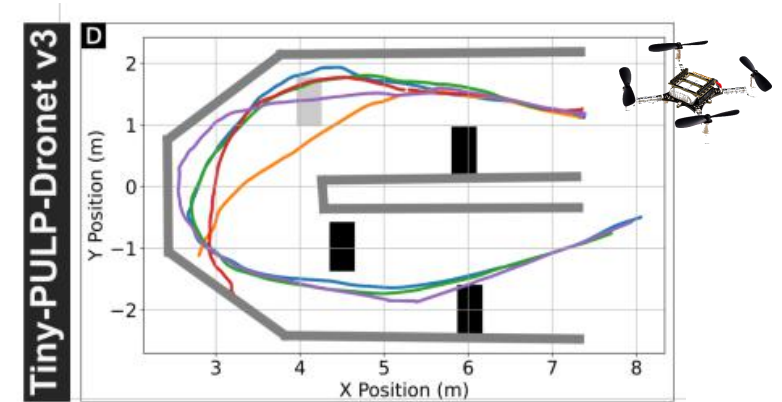
<https://arxiv.org/ftp/arxiv/papers/2311/2311.01057.pdf>

## DSP/NN Mixed-Precision – Denoiser



<https://arxiv.org/pdf/2210.07692>

## Nano drones autonomous navigation



<https://arxiv.org/pdf/2407.12675>

# TinyML: best performance HW and SW

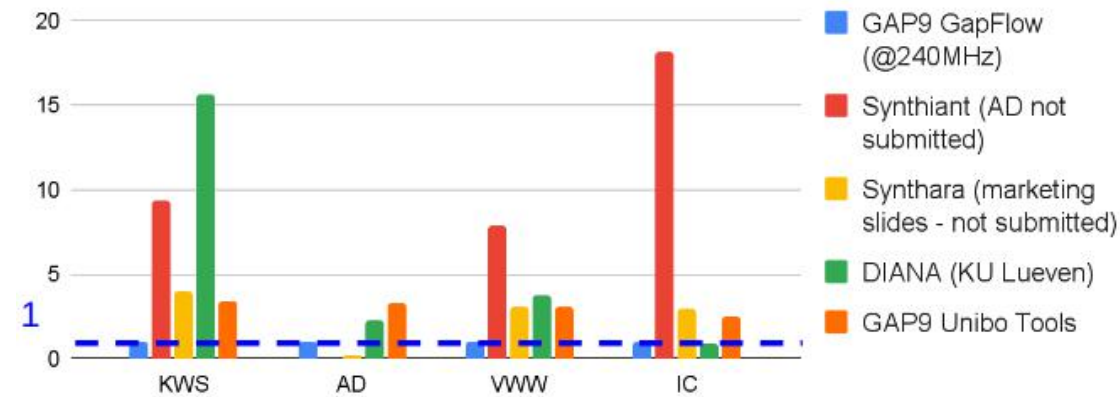
## Latency:

- 7-18x faster than second best submitted result (Synthiant)
- Up to 3x faster than industry leader yet to be available (Synthara)
- Up to 15x faster than latest research chips (DIANA)
- 2-4x faster than other publicly available SW libraries/tools with same HW (GAP9)

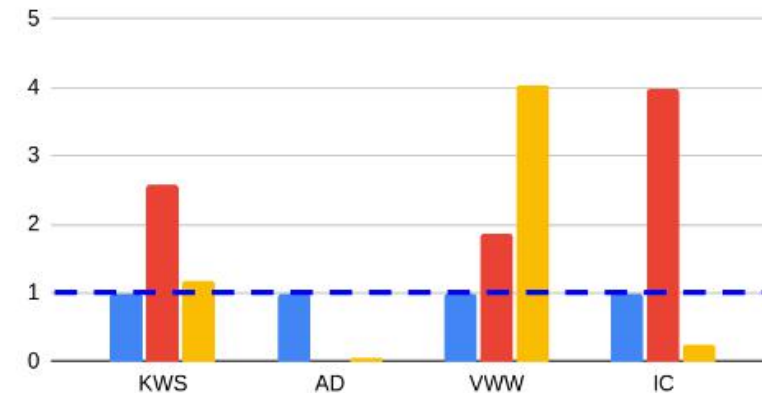
## Energy:

- 2-5x lower energy than second best submitted result
- Up to 4x better than not yet available HW
- \*DIANA and GAP9 Unibo data not available for energy

**Latency normalized to GAP9 GAPflow**



**Energy normalized to GAP9 GAPflow**



# Workshop

# Exercise:

1. Create graph in nntool to compute a matrix multiplication ( $C = A*B$ ) using int8 quantization and profile it in different scenarios:
  - a) A (32x64) and B (64x128) variable
  - b) A (32x64) variable and B (64x128) constant - is it better than a)? Why? Can you improve it more?
2. Add a resize node in front of the mobilenet v2 we deployed (DO NOT use the PATCH trick)
  - a) Force the input to be stored in L3 RAM
3. Deploy a yolox model like we did for the mobilenet v2. Can all the optimizations we did be applied straight-forward? (maybe you have to change something :)

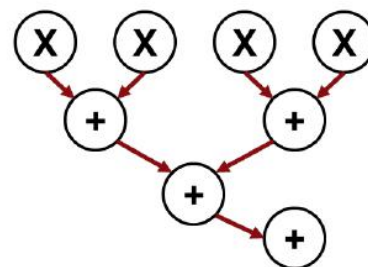
# Convolution HW acceleration: NE16

# NN Digital Accelerators

- Two main families:

- **Convolvers:**

- Filter specific (difficult to generalize)
    - Maximal data reusage
    - Work well in Depth-wise convolutions
    - Examples: NE16



Implement the spatial filter as an **adder tree**

- **MatMul Accelerators:**

- Adapts well to any type of filter sizes
    - Does not work in Depth-wise
    - Requires Im2Col
    - Most commonly used approach
    - Examples: Google TPU, NVDLA

- **More Options available:**

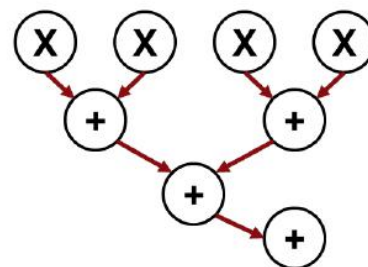
- **Combining** ideas of the 2 families
    - Exploit **serialization**: e.g. bit-serial multipliers
    - Exploit **sparsity**

# NN Digital Accelerators

- Two main families:

- **Convolvers:**

- Filter specific (difficult to generalize)
    - Maximal data reusage
    - Work well in Depth-wise convolutions
    - Examples: NE16



Implement the spatial filter as an **adder tree**

- **MatMul Accelerators:**

- Adapts well to any type of filter sizes
    - Does not work in Depth-wise
    - Requires Im2Col
    - Most commonly used approach
    - Examples: Google TPU, NVDLA

- **More Options available:**

- **Combining** ideas of the 2 families
    - Exploit **serialization**: e.g. bit-serial multipliers
    - Exploit **sparsity**

**NE16**

# Why bit serialization?

$$\begin{array}{r}
 01011101 \\
 11001000 \\
 \hline
 \end{array}
 \times$$

$$1100100011001000$$

**One-step** solution:

- more **power efficient**
- more area (**more complex**)
- not scalable (less bit, same power)  
e.g. 4x8bits = 8x8bits

**Bit-Serial** solution:

- Each cycle compute 1bit product
- less area (**simpler**)
- less bits, less cycles (energy)  
e.g. 4x8bits ~ 1/2 time of 8x8bits

$$\begin{array}{r}
 01011101 \\
 11001000 \\
 \hline
 \end{array}$$

$$00000000$$

**shift & add** 01011101

$$010111010$$

**shift & add** 00000000

⋮

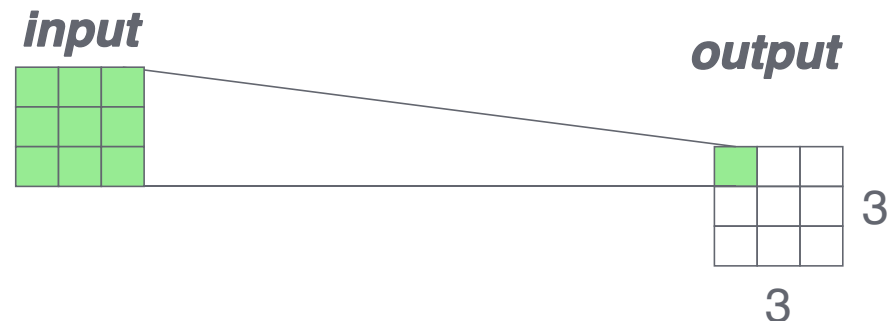
x8





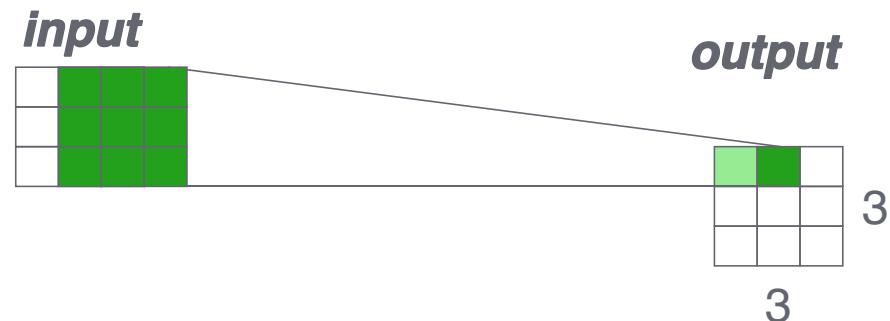
## NE16 - 3x3 Convolver

- Why 3x3? Widely adopted in many vision NN (SqueezeNet, VGG, ...)
- **Design choices:**
  - **3x3 filter** computed in parallel --> each adder tree has 9 filter elements
  - **3x3 output pixels** computed in **parallel** --> 9 adder trees
  - **16 input channels** computed in **parallel** --> each adder tree has 16 channels



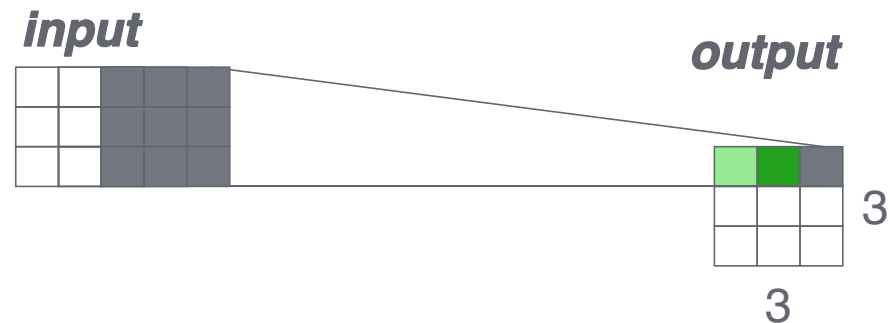
## NE16 - 3x3 Convolver

- Why 3x3? Widely adopted in many vision NN (SqueezeNet, VGG, ...)
- **Design choices:**
  - **3x3 filter** computed in parallel --> each adder tree has 9 filter elements
  - **3x3 output pixels** computed in **parallel** --> 9 adder trees
  - **16 input channels** computed in **parallel** --> each adder tree has 16 channels



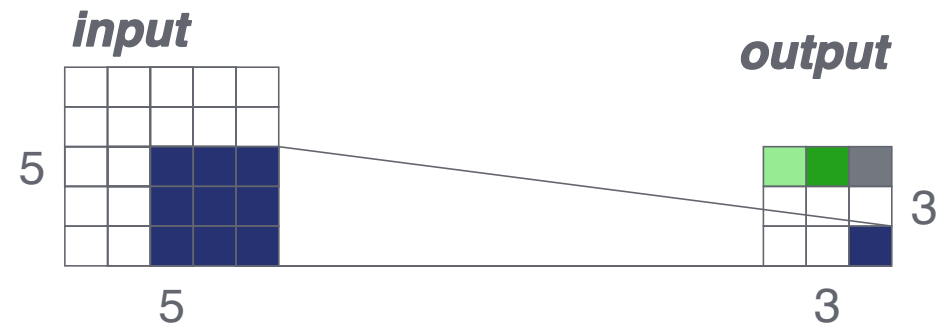
## NE16 - 3x3 Convolver

- Why 3x3? Widely adopted in many vision NN (SqueezeNet, VGG, ...)
- **Design choices:**
  - **3x3 filter** computed in parallel --> each adder tree has 9 filter elements
  - **3x3 output pixels** computed in **parallel** --> 9 adder trees
  - **16 input channels** computed in **parallel** --> each adder tree has 16 channels



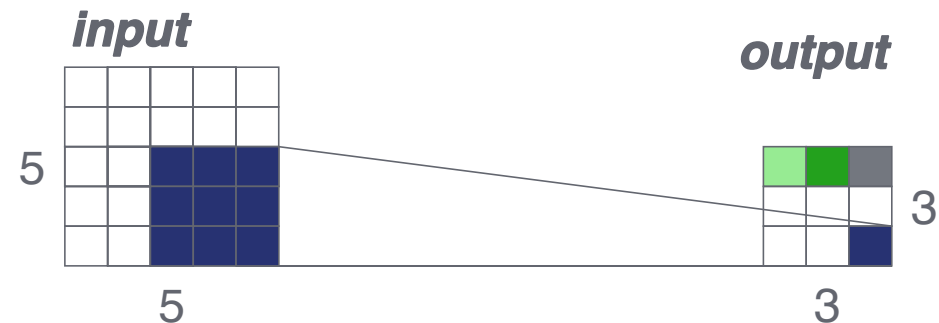
## NE16 - 3x3 Convolver

- Why 3x3? Widely adopted in many vision NN (SqueezeNet, VGG, ...)
- **Design choices:**
  - **3x3 filter** computed in parallel --> each adder tree has 9 filter elements
  - **3x3 output pixels** computed in **parallel** --> 9 adder trees
  - **16 input channels** computed in **parallel** --> each adder tree has 16 channels



## NE16 - 3x3 Convolver

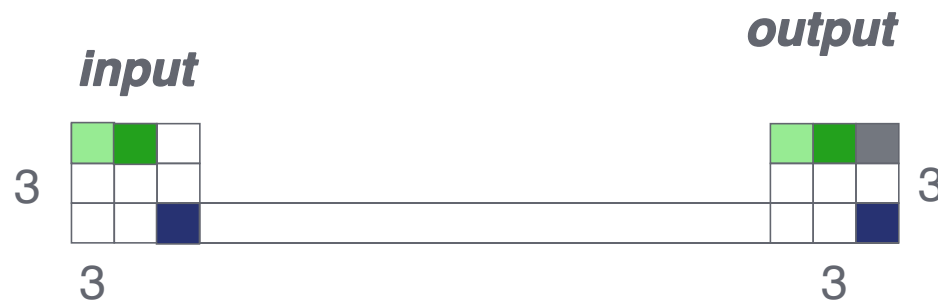
- Why 3x3? Widely adopted in many vision NN (SqueezeNet, VGG, ...)
- **Design choices:**
  - **3x3 filter** computed in parallel --> each adder tree has 9 filter elements
  - **3x3 output pixels** computed in **parallel** --> 9 adder trees
  - **16 input channels** computed in **parallel** --> each adder tree has 16 channels



- **5x5 input buffer**
- **32 output channels before reiterating**

## NE16 - 3x3 Convolver (ONLY?)

- How does it do 1x1 (aka MatMul)?
  - **1x1 filter** computed in parallel --> each adder tree has 1 filter elements
  - **Parallelize on the weight bits (no more bit serial here!!!)** --> 2-8 bits in parallel
  - **3x3 output pixels** computed in **parallel** --> 9 adder trees
  - **16 input channels** computed in **parallel** --> each adder tree has 16 channels



- **5x5 input buffer**
- **32 output channels before reiterating**

# Theory is cool, but reality?

	<b>3x3</b>	<b>1x1</b>
<b>Theory</b>	162 MAC/Cyc	144 MAC/Cyc
<b>Reality</b>	~130 MAC/Cyc (~80%)	~60 MAC/Cyc (<50%)



## Why is 1x1 so bad in reality?

- In both modes before starting the accumulation we need to load the input buffer !!!!

	<i>Input buffer size</i>	<i>Number of MACs for the input buffer</i>	<i>MAC/Load ratio</i>
<b>3x3</b>	5x5x16	$32 \times 3 \times 3 \times 16 \times 3 \times 3 = 41472$	103.7
<b>1x1</b>	3x3x16	$32 \times 3 \times 3 \times 16 = 4608$	32.0

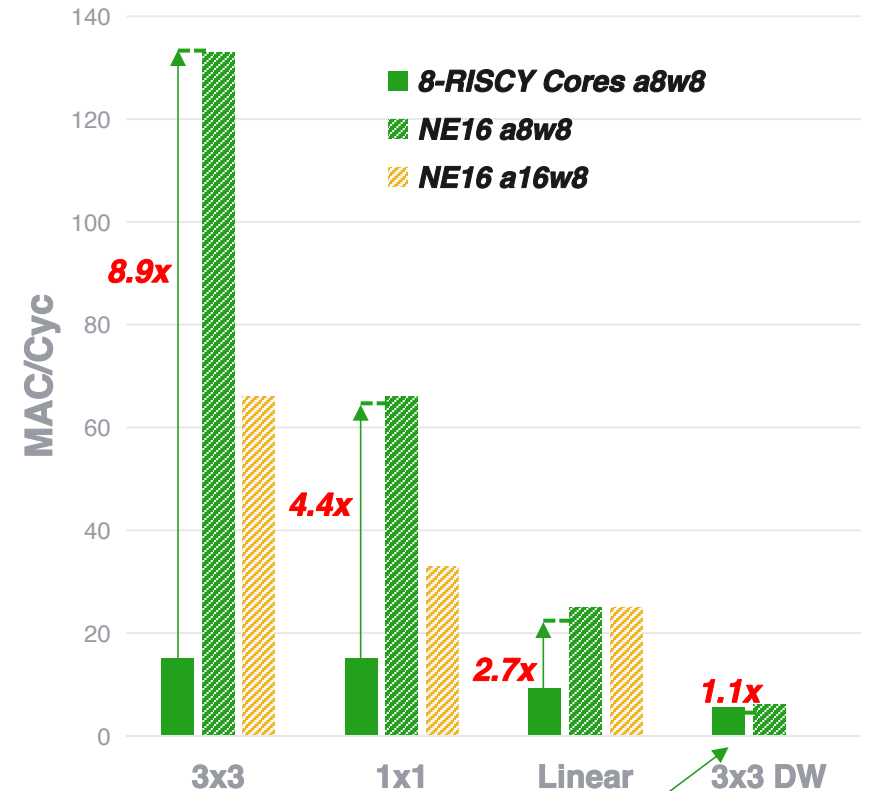
1x1 mode has lower MAC/Load ratio --> it could go faster but it is bounded by the input buffer load



# NE16HW accelerator – heterogeneous performance

- Operating modes:
  - **3x3** w/ pad (w/ stride 1x1 or 2x2)
  - **1x1** (MatMul)
  - **Linear** (MatVector multiply)
  - **DepthWise 3x3** w/ pad (w/ stride 1x1 or 2x2)
- NE SW Library:
  - We can combine parallel SW (8 cores) and NE16 to achieve the best performance in many scenarios, e.g. Im2Col + NE16-1x1

**NE16 Real Use Case Performance**

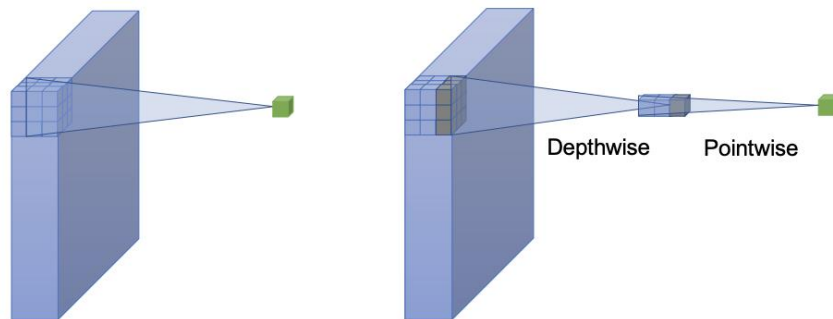


No perf gain but allows to keep the HWC order

# Architecture design still matters

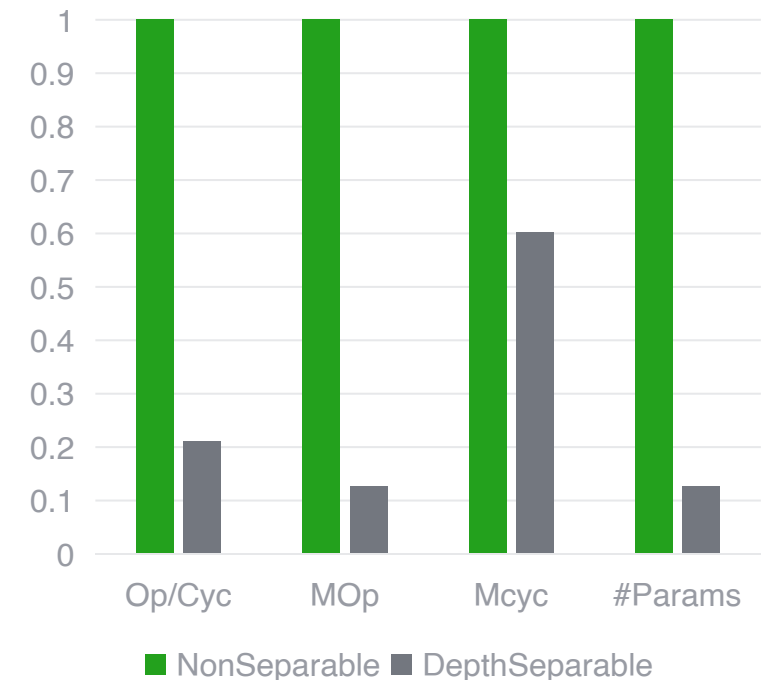
Even with dedicated HW, NN architecture can help AI on the edge, e.g.:

- **3x3 Non-Separable** Convolution:
  - NE16-friendly --> **120MAC/Cyc**
- **3x3 DepthSeparable** Convolution: **3x3 DW** + **1x1 PW**
  - non NE16-friendly --> **6MAC/Cyc** + **60MAC/Cyc** ~ **20MAC/Cyc**



[Depthwise Convolution is All You Need for Learning Multiple Visual Domains]

Standard vs DepthSeparable 3x3 Convolution on NE16



# Hands-on