#### A 12.4 TOPS/W @ 136GOPS AI-IoT System-on-Chip with 16 RISC-V, 2-to-8b Precision-Scalable DNN Acceleration and 30%-Boost Adaptive Body Biasing

<u>Francesco Conti<sup>1</sup></u>, Davide Rossi<sup>1</sup>, Gianna Paulin<sup>2</sup>, Angelo Garofalo<sup>1</sup>, Alfio Di Mauro<sup>2</sup>, Georg Rutishauer<sup>2</sup>, Gianmarco Ottavi<sup>1</sup>, Manuel Eggimann<sup>2</sup>, Hayate Okuhara<sup>1</sup>, Vincent Huard<sup>3</sup>, Olivier Montfort<sup>3</sup>, Lionel Jure<sup>3</sup>, Nils Exibard<sup>3</sup>, Pascal Gouedo<sup>3</sup>, Mathieu Louvat<sup>3</sup>, Emmanuel Botte<sup>3</sup>, Luca Benini<sup>1,2</sup>



<sup>1</sup>University of Bologna, <sup>2</sup>ETH Zürich, <sup>3</sup>Dolphin Design







# Introduction and Motivation

- Emerging application areas for Al-enabled IoT
  - Personalized Healthcare
  - Augmented Reality
  - Nano-Robotics
- Challenges
  - High computational demand from DNNs, other algorithms
  - Diverse computational patterns and requirements
  - Opportunities
    - Bit-precision tolerance
    - Accelerable workloads



# Introduction and Motivation

- Emerging application areas for Al-enabled IoT
  - Personalized Healthcare
  - Augmented Reality
  - Nano-Robotics
- Challenges
  - High computational demand from DNNs, other algorithms
  - Diverse computational patterns and requirements
  - Opportunities
    - Bit-precision tolerance
    - Accelerable workloads



# Introduction and Motivation

- Emerging application areas for Al-enabled IoT
  - Personalized Healthcare
  - Augmented Reality
  - Nano-Robotics

#### Challenges

- High computational demand from DNNs, other algorithms
- Diverse computational patterns and requirements
- Opportunities
  - Bit-precision tolerance
  - Accelerable workloads



- Programmable RISC-V based System-on-Chip combining
  - 1. RISC-V MCU with 1MB of on-board SRAM
    - 32-bit RV32IMCFXpulp
    - 1MB L2 SRAM
    - Standard peripherals (SPI, I2C, UART, ...)
    - Off-chip \*HyperRAM™
       DRAM/FLASH
    - Autonomous I/O DMA
    - 3 Freq-Locked Loops



#### \*https://www.cypress.com/products/hyperram-memory

- Programmable RISC-V based System-on-Chip combining
  - 1. RISC-V MCU with 1MB of on-board SRAM
  - 2. Cluster of 16 RISC-V 2-32b DSP cores
    - 16x RV32IMCF<u>XpulpNN</u> Period
    - 2b/4b + fused MAC&LOAD
    - 8x FP32/FP16/BF16 FPUs
    - 128 KB of L1 TCDM (0-wait-state, 400Gb/s @ 420MHz aggr. bandwidth)
    - DMA w/ 64b L2 $\rightarrow$ L1 + 64b L1 $\rightarrow$ L2  $\rightarrow$ C



- Programmable RISC-V based System-on-Chip combining
  - 1. RISC-V MCU with 1MB of on-board SRAM
  - 2. Cluster of 16 RISC-V 2-32b DSP cores
  - 3. 2-8b Reconfigurable
     Binary Engine for 3x3,
     1x1 DNN kernels
    - asymmetric 2-8b inputs (*I-b*), weights (*W-b*), outputs (*O-b*) composing 1x1-b products
    - up to 10<sup>4</sup> 1x1b-MAC/cycle



- **Programmable RISC-V based** System-on-Chip combining
  - 1. **RISC-V MCU** with 1MB of on-board SRAM
  - 2. Cluster of 16 RISC-V 2-32b DSP cores
  - 2-8b Reconfigurable 3. **Binary Engine for 3x3,** 1x1 DNN kernels
  - Flash **Adaptive Body Biasing** 4. with on-the-fly control
    - OCMs added to 1% most timing-critical endpoints
    - Adaptively boost freq or efficiency



22.1: A 12.4TOPS/W @ 136GOPS AI-IoT System-on-Chip with 16 RISC-V 2-to-8b Precision-Scalable DNN Acceleration and 30%-Boost Adaptive Body Biasing

Bank

#30

Bank

#31

#15



- **Programmable RISC-V based** System-on-Chip combining
  - **RISC-V MCU** with 1MB of 1. on-board SRAM
  - 2. Cluster of 16 RISC-V 2-32b DSP cores
  - 2-8b Reconfigurable 3. **Binary Engine for 3x3,** 1x1 DNN kernels
  - **Adaptive Body Biasing** 4. with on-the-fly control
- **Prototype implemented in GF 22FDX** 
  - > flip-well LVT & SLVT cells, 2.43mm<sup>2</sup> for CLUSTER

1.94 mm



3.78 mm

#### **Baseline RISC-V Core**



#### **RISC-V MAC&LD core architecture**



Xpulp Dot-Product achieves up to 2.3 MAC/cycle/core @ 8-bit

lp.setup			
p.lw w0, 4(	(WO!)		
p.lw w1, 4(	(w1!)		
p.lw w2, 4(	(W2!)		
p.lw w3, 4(	(W3!)		
p.lw x1, 4(	(X0!)		
p.lw x2, 4(	(X1!)		
pv.sdotsp.b	acc1,	w0,	x0
pv.sdotsp.b	acc2,	w0,	x1
pv.sdotsp.b	acc3,	w1,	<b>x0</b>
pv.sdotsp.b	acc4,	w1,	x1
pv.sdotsp.b	acc5,	w2,	<b>x0</b>
pv.sdotsp.b	acc6,	w2,	x1
pv.sdotsp.b	acc7,	w3,	<b>x0</b>
pv.sdotsp.b	acc8,	w3,	x1
end	<u>2.29 M</u>	<u>AC/c</u>	<u>ycle/core</u>

### **RISC-V MAC&LD core architecture**



*XpulpNN* adds new Dot-Product instructions that exploit a dedicated NN-RF

- 1. to boost weight/input reuse
- 2. to overlap loads and dotproduct

lp.setup		
pv.nnsdotup.h zero,	x1,	9
pv.nnsdotsp.b acc1, N	<b>ν2</b> ,	0
pv.nnsdotsp.b acc2, N	<b>ν4</b> ,	2
pv.nnsdotsp.b acc3, N	<b>ν3</b> ,	4
pv.nnsdotsp.b acc4, x	x1,	14
pv.nnsdotsp.b acc5, N	<b>ν2</b> ,	17
pv.nnsdotsp.b acc6, N	<b>ν4</b> ,	19
pv.nnsdotsp.b acc7, N	<b>ν3</b> ,	21
pv.nnsdotsp.b acc8, N	w1,	23
end <u>3.56 MAC/cycle (+55</u>	<mark>% b</mark>	<mark>oost</mark> )



- Partially bit-serial dataflow
- "Bit" dimension is decomposed, in inputs, weights, outputs



- DNN layer operating at *WxI*-b decomposed in *WxI* 1x1b MAC ops
  - RBE tensors use a special layout in memory to expose 1x1b-ops
- RBE is integrated in cluster with
  - 288b LD/ST initiator port toward CLUSTER interconnect
  - 32b config port + "end of job" event



#### RBE is 652 kGE divided in

- Register File + Control (6%)
- 288b Load/Store Unit (2%)
- Datapath (92%)

- Datapath comprise 9 Cores
  - 1 Core = receptive field of one output in space across 32 Kout
  - Stationary input statically multicasted to each Core receptive field
  - Output stored in 32x 32b accumulators per Core + Quantized
  - Weights renovated each cycle and broadcasted on all Cores



9 Blocks per Core

- Each with 4 <u>Bin</u>ary <u>Conv</u>olvers with 32 1x1b "multipliers" (i.e., AND gates)
- Results from all BinConvs are scaled, cumulated with Adder Trees to form a complete contribution

BinConv

BinConv



### **Frequency & Power**

- **MARSELLUS** is operational within 0.5-0.8V voltage range
- Frequency range with no Vbb 100 MHz @ 0.5V → 420 MHz @ 0.8V
- Frequency [MHz] **CLUSTER** power ranges between 12.8mW @ 100MHz, 0.5V → 123mW @ 420MHz, 0.8V







### **Efficiency Sweep**

SW Workload: MatMul no sparsity



Baseline efficiency with 16 cores not using M&L is 130 Gop/s/W @ 0.8V → 290 Gop/s/W @ 0.5V

#### 

## **Efficiency Sweep**



- Baseline efficiency with 16 cores not using M&L is 130 Gop/s/W @ 0.8V
   → 290 Gop/s/W @ 0.5V
- M&L and quantization down to 2b bring up to 12.8x improvement in efficiency

#### → MMUL 8b → MMUL M&L 8b → MMUL M&L 4b → MMUL M&L 2b

SW Workload: MatMul

# **Efficiency Sweep**



Baseline efficiency with 16 cores not using M&L is 130 Gop/s/W @ 0.8V → 290 Gop/s/W @ 0.5V M&L and quantization down to 2b bring up to 12.8x improvement in efficiency Employing RBE on

Employing RBE on CONV3x3 kernels efficiency can be improved by a further 7.6x, up to 12.4Top/s/W

© 2023 IEEE International Solid-State Circuits Conference HW Workload: Conv3x3 32x16x3x3

- On-Chip Monitors detect preerror conditions
  - 1% most critical endpoints
  - shadow delayed FF compared with regular FF → difference = pre-error



- On-Chip Monitors detect preerror conditions
  - 1% most critical endpoints
  - shadow delayed FF compared with regular FF → difference = pre-error
- When OCM detect pre-errors,
   ABB generators adaptively tweak body biasing



- On-Chip Monitors detect preerror conditions
  - 1% most critical endpoints
  - shadow delayed FF compared with regular FF → difference = pre-error
- When OCM detect pre-errors, ABB generators adaptively tweak body biasing
  - Either boost frequency up to
     12% (420 → 470MHz)



- On-Chip Monitors detect preerror conditions
  - 1% most critical endpoints
  - shadow delayed FF compared with regular FF → difference = pre-error
- When OCM detect pre-errors,
   ABB generators adaptively
   tweak body biasing
  - Either boost frequency up to 12% (420 → 470MHz)
  - Or keep a fixed frequency target (e.g., 400 MHz) while dropping Vdd (0.8 → 0.65V) → +30% efficiency

#### ABB-enabled V<sub>DD</sub> down-scaling (f=400MHz)



22.1: A 12.4TOPS/W @ 136GOPS AI-IoT System-on-Chip with 16 RISC-V 2-to-8b Precision-Scalable DNN Acceleration and 30%-Boost Adaptive Body Biasing

# **Efficiency Techniques in MARSELLUS**



© 2023 IEEE International Solid-State Circuits Conference 22.1: A 12.4TOPS/W @ 136GOPS AI-IoT System-on-Chip with 16 RISC-V 2-to-8b Precision-Scalable DNN Acceleration and 30%-Boost Adaptive Body Biasing

# **Efficiency Techniques in MARSELLUS**



© 2023 IEEE International Solid-State Circuits Conference

<sup>22.1:</sup> A 12.4TOPS/W @ 136GOPS AI-IoT System-on-Chip with 16 RISC-V 2-to-8b Precision-Scalable DNN Acceleration and 30%-Boost Adaptive Body Biasing

# **Efficiency Techniques in MARSELLUS**



<sup>22.1:</sup> A 12.4TOPS/W @ 136GOPS AI-IoT System-on-Chip with 16 RISC-V 2-to-8b Precision-Scalable DNN Acceleration and 30%-Boost Adaptive Body Biasing

#### Full network: ResNet-20/CIFAR-10 on RBE



#### Full network: ResNet-20/CIFAR-10 on RBE



- Acc/Energy for CONV layers (no ADD)
- In, Out, ADD quantized to 8-bits
- Can run at up to 4000fps @ 0.65V+ABB
- 20 uJ/frame
- Or at 1000fps @ 0.5V
  - 12 uJ/frame

ACCURACY [%] ENERGY [uJ] LATENCY [us]

#### State-of-the-Art Comparison

	<b>VEGA</b> [1]	SAMURAI [2]	DIANA [3]	<b>QNAP [4]</b>	MARSELLUS
	ISSCC'21	VLSIC'20	ISSCC'22	<b>ISSCC'21</b>	(this work)
Technology	22nm FDX	28nm FD-SOI	22nm FDX + AIMC	28nm	22nm FDX
Die Area	10 mm <sup>2</sup>	4.5 mm <sup>2</sup>	10.24 mm <sup>2</sup>	1.9 mm <sup>2</sup>	18.7 mm <sup>2</sup> (cluster 1.9 mm <sup>2</sup> )
Cores	10x RV32IMCFXpulp + HWCE	1x RV32IMCFXpulp + Digital Accel.	1x RV32CIMFXpulp + Digital Accel. + AIMC	Digital Accel.	1x RV32IMCFXpulp + 16x RV32IMCFXpulpnn + RBE
Max Frequency	450 MHz	350 MHz	320 MHz	470 MHz	420 MHz
Power range	1.7 uW - 49.4 mW	6.4 uW - 96 mW	10-129 mW (digital)	19.4-131 mW	12.8 mW - 123 mW
Best SW (INT) Perf	15.6 GOPS (8 RISC-V)	1.5 GOPS (1 RISC-V)	-	-	90 GOPS (16 RISC-V M&L 2x2b, 0.8V+ABB)
Best SW (INT) Eff	614 GOPS/W @ 7.6 GOPS (8 RISC-V)	230 GOPS/W @ 110 MOPS (1 RISC-V)	-	-	1.66 TOPS/W @ 19 GOPS (16 RISC-V M&L 2x2b)
Best HW-Accel Perf	32.2 GOPS (HWCE)	36 GOPS (Digital Accel.)	180 GOPS (Digital), 29.5 TOPS (AIMC)	140 GOPS (Digital Accel.)	637 GOPS (RBE 2x2b, 0.8V+ABB)
Best HW-Accel Eff	1.3 TOPS/W @ 15.6 GOPS (HWCE)	1.3 TOPS/W @ 2.8 GOPS (Digital Accel.)	4.1 TOPS/W (Digital), 600 TOPS/W (AIMC)	12.1 TOPS/W @ 140 GOPS (Digital Accel.)	12.4 TOPS/W @ 136 GOPS (RBE 2x2b)

#### **State-of-the-Art Comparison**

	VEGA [1]	SAMURAI [2]	DIANA [3]	<b>QNAP [4]</b>	MARSELLUS
	ISSCC'21	VLSIC'20	ISSCC'22	ISSCC'21	(this work)
Technology	22nm FDX	28nm FD-SOI	22nm FDX + AIMC	28nm	22nm FDX
Die Area	10 mm <sup>2</sup>	4.5 mm <sup>2</sup>	10.24 mm <sup>2</sup>	1.9 mm <sup>2</sup>	18.7 mm <sup>2</sup> (cluster 1.9 mm <sup>2</sup> )
Cores	10x RV32IMCFXpulp + HWCE	1x RV32IMCFXpulp + Digital Accel.	1x RV32CIMFXpulp + Digital Accel. + AIMC	Digital Accel.	1x RV32IMCFXpulp + 16x RV32IMCFXpulpnn + RBE
Max Frequency	450 MHz	350 MHz	320 MHz	470 MHz	420 MHz
Power range	1.7 uW 49.4 mW	6.4 uW - 96 mW	10-129 mW (digital)	19.4-131 mW	12.0 mW 123 mW
Best SW (INT) Per	15.6 GOPS (8 RISC-V)	1.5 GOPS (1 RISC-V)	Machine Learni	ng SW	90 GOPS (16 RISC-V M&L 2x2b, 0.8V+ABB)
Best SW (INT) Eff	614 GOPS/W @ 7.6 GOPS (8 RISC-V)	230 GOPS/W @ 110 MOPS (1 RISC-V)	of 5.8x & 2.7 compared to V	7boost 7x EGA	1.66 TOPS/W @ 19 GOPS (16 RISC-V M&L 2x2b)
Best HW-Accel Per	(HWCE)	36 GOPS (Digital Accel.)	180 GOPS (Digital), 29.5 TOPS (AIMC)	140 GOPS (Digital Accel.)	637 GOPS (RBE 2x2b, 0.8V+ABB)
Best HW-Accel Eff	1.3 TOPS/W @ 15.6 GOPS (HWCE)	1.3 TOPS/W @ 2.8 GOPS (Digital Accel.)	4.1 TOPS/W (Digital), 600 TOPS/W (AIMC)	12.1 TOPS/W @ 140 GOPS (Digital Accel.)	12.4 TOPS/W @ 136 GOPS (RBE 2x2b)

#### **State-of-the-Art Comparison**

	<b>VEGA</b> [1]	SAMURAI [2]	DIANA [3]	<b>QNAP [4]</b>	MARSELLUS
	ISSCC'21	VLSIC'20	<b>ISSCC'22</b>	<b>ISSCC'21</b>	(this work)
Technology	22nm FDX	28nm FD-SOI	22nm FDX + AIMC	28nm	22nm FDX
Die Area	10 mm <sup>2</sup>	4.5 mm <sup>2</sup>	10.24 mm <sup>2</sup>	1.9 mm <sup>2</sup>	18.7 mm <sup>2</sup> (cluster 1.9 mm <sup>2</sup> )
Cores	10x RV32IMCFXpulp + HWCE	1x RV32IMCFXpulp + Digital Accel.	1x RV32CIMFXpulp + Digital Accel. + AIMC	Digital Accel.	1x RV32IMCFXpulp + 16x RV32IMCFXpulpnn + RBE
Max Frequency	450 MHz	350 MHz	320 MHz	470 MHz	420 MHz
Power range	1.7 uW - 49.4 mW	6.4 uW - 96 mW	10-129 mW (digital)	19.4-131 mW	12.8 mW - 123 mW
Best SW (INT) Perf	15.6 GOPS (8 RISC-V)	1.5 GOPS (1 RISC-V)	- I Per	DNN Acceleration	90 GOPS 1 (16 RISC-V M&L 2x2b, ost 0.8V+ABB)
Best SW (INT) Eff	614 GOPS/W @ 7.6 GOPS (8 RISC-V)	230 GOPS/W @ 110 MOPS (1 RISC-V)	- compa	3.5x & 3.0x ared to DIANA (D	1.56 TOPS/W igital) @ 19 GOPS 1/16 PISC-V M&L 2x2b)
Best HW-Accel Perf	32.2 GOPS (HWCE)	36 GOPS (Digital Accel.)	180 GOPS (Digital), 29.5 TOPS (AIMC)	140 GOPS (Digital Accel.)	637 GOPS (RBE 2x2b, 0.8V+ABB)
Best HW-Accel Eff	1.3 TOPS/W @ 15.6 GOPS (HWCE)	1.3 TOPS/W @ 2.8 GOPS (Digital Accel.)	4.1 TOPS/W (Digital), 600 TOPS/W (AIMC)	12.1 TOPS/W @ 140 GOPS (Digital Accel.)	12.4 TOPS/W @ 136 GOPS (RBE 2x2b)

#### **State-of-the-Art Comparison (2)**

	VEGA [1] ISSCC'21	SAMURAI [2] VLSIC'20	DIANA [3] ISSCC'22	QNAP [4] ISSCC'21	MARSELLUS (this work)
Best ResNet-20 Eff (CIFAR-10)	-	-	14.4 TOPS/W (AIMC)	-	6.38 TOPS/W (RBE mixed)
ResNet-20 Latency @ Best Eff	-	-	1.26 ms	-	1.05 ms
Best ResNet-18 Eff (ImageNet)	-	-	19 TOPS/W	12.62 TOPS/W (zero-skipping)	5.83 TOPS/W (RBE 4x4b)
ResNet-18 Latency @ Best Eff	-	-	6.15ms	24.8ms	48ms

### **State-of-the-Art Comparison (2)**

	VEGA [1] ISSCC'21	SAMURAI [2] VLSIC'20	DIANA [3] ISSCC'22	QNAP [4] ISSCC'21	MARSELLUS (this work)
Best ResNet-20 Eff (CIFAR-10)	-	-	14.4 TOPS/W (AIMC)	-	6.38 TOPS/W (RBE mixed)
ResNet-20 Latency @ Best Eff	-	-	1.26 ms	-	1.05 ms
Best ResNet-18 Eff (ImageNet)	-	-	19 TOPS/W	12.62 TOPS/W (zero-skipping)	5.83 TOP S/W (RBE 4x4b)
ResNet-18 Latency @ Best Eff	-	-	6.15ms	24.8ms	48ms

#### ResNet-20/CIFAR10: digital MARSELLUS vs mixed-signal AIMC-based accelerator DIANA

- ~40x lower performance and efficiency at peak, but just 2.3x worse efficiency at 20% better latency
- Marginal energy advantage might be further neglected in full system

### State-of-the-Art Comparison (2)

	VEGA [1] ISSCC'21	SAMURAI [2] VLSIC'20	DIANA [3] ISSCC'22	QNAP [4] ISSCC'21	MARSELLUS (this work)
Best ResNet-20 Eff (CIFAR-10)	-	-	14.4 TOPS/W (AIMC)	-	6.38 TOPS/W (RBE mixed)
ResNet-20 Latency @ Best Eff	-	-	1.26 ms	-	1.05 ms
Best ResNet-18 Eff (ImageNet)	-	-	19 TOPS/W	12.62 TOPS/W (zero-skipping)	5.83 TOPS/W (RBE 4x4b)
ResNet-18 Latency @ Best Eff	-	-	6.15ms	24.8ms	48ms

- ResNet-20/CIFAR10: digital MARSELLUS vs mixed-signal AIMC-based accelerator DIANA
  - ~40x lower performance and efficiency at peak, but just 2.3x worse efficiency at 20% better latency
  - Marginal energy advantage might be further neglected in full system
- MARSELLUS shows competitive acceleration even compared with AIMC for TinyML and AIloT endpoint devices
  - Up to 5.83 TOPS/W on 4-bit ImageNet @ 20fps without AIMC, zero-skipping, sparsity

### Conclusion

#### 4.95 mm



#### MARSELLUS combines

- state-of-the-art parallel and heterogeneous acceleration
- aggressive performance and voltage scalability thanks to ABB

#### MARSELLUS achieves

1000x better performance and efficiency than commercial AI-IoT microcontrollers with similar power budget (~130mW)

# Thanks for the attention. Questions?

#### References

- [1] D. Rossi et al., "A 1.3TOPS/W @ 32GOPS Fully Integrated 10-Core SoC for IoT End- Nodes with 1.7µW Cognitive Wake-Up From MRAM-Based State-Retentive Sleep Mode," *ISSCC*, pp. 60-62, 2021.
- [2] I. Miro-Panades et al., "SamurAI: A 1.7MOPS-36GOPS Adaptive Versatile IoT Node with 15,000× Peak-to-Idle Power Reduction, 207ns Wake-Up Time and 1.3TOPS/W ML Efficiency," *IEEE Symp. VLSI Circuits*, 2020.
- [3] K. Ueyoshi et al., "DIANA: An End-to-End Energy-Efficient Digital and ANAlog Hybrid Neural Network SoC," ISSCC, pp. 256-257, 2022.
- [4] H. Mo et al., "A 28nm 12.1TOPS/W Dual-Mode CNN Processor Using Effective- Weight-Based Convolution and Error-Compensation-Based Prediction," *ISSCC*, pp. 146-148, 2021.
- [5] Z. Dong et al., "HAWQ: Hessian AWare Quantization of Neural Networks With Mixed-Precision," *ICCV*, pp. 293-302, 2019.

# Appendix: RBE multiple-precision performance



© 2023 IEEE International Solid-State Circuits Conference 22.1: A 12.4TOPS/W @ 136GOPS AI-IoT System-on-Chip with 16 RISC-V 2-to-8b Precision-Scalable DNN Acceleration and 30%-Boost Adaptive Body Biasing

### **Appendix: RBE Execution flow**



22.1: A 12.4TOPS/W @ 136GOPS AI-IoT System-on-Chip with 16 RISC-V 2-to-8b Precision-Scalable DNN Acceleration and 30%-Boost Adaptive Body Biasing

### **Appendix: RBE CONV mapping**



© 2023 IEEE International Solid-State Circuits Conference 22.1: A 12.41 UPS/W @ 136GOPS Al-IoT System-on-Chip with 16 RISC-V 2-to-8b Precision-Scalable DNN Acceleration and 30%-Boost Adaptive Body Biasing

### **Appendix: RBE Dataflow loop nests**

```
3x3 CONV
                                                                                                     1x1 CONV
for(h in=0 to 5) :
                                                        for(h in=0 to 3) :
                                                         for(bit in=0 to I) :
 for(bit in=0 to I) :
                                                            par for(w in=0 to 3) : // due to Activation Layout
   par_for(w_in=0 to 5) : // due to Activation Layout
                                                            par for(k in=0 to 32) : // due to Activation Layout
   par_for(k_in=0 to 32) : // due to Activation Layout
                                                              load(in buffer[h in, bit in, w in, k in])
     load(in_buffer[h_in, bit_in, w_in, k_in])
                                                             inp = map(in buffer)
      inp = map(in buffer)
                                                        for(k out=0 to 32) :
for(k out=0 to 32) :
                                                         par for(h out=0 to 3) : // over Cores
  for(bit wgt=0 to W) :
                                                         par for(w out=0 to 3) : // over Cores
   par for(h out=0 to 3) : // over Cores
                                                         par for(bit wgt=0 to W) : // over rows of Blocks
   par for(w out=0 to 3) : // over Cores
                                                         par for(bit in=0 to I) : // over BinConvs
   par for(f=0 to 9) : // over rows of Blocks
                                                         par_for(k_in=0 to 32) : // inside BinConv
   par for(bit in=0 to I) : // over BinConvs
                                                            load(wgt[k out, k in tiles, bit wgt, f, k in])
   par for(k in=0 to 32) : // inside BinConv
                                                            acc[h_out, w_out, k_out] += dot(wgt, inp)
     load(wgt[k_out, k_in_tiles, bit_wgt, f, k_in])
      acc[h out, w out, k out] += dot(wgt, inp)
```

LOAD

COMPUTE