
A 1.15 TOPS/W, 16-Cores Parallel Ultra-Low Power Cluster with 2b-to-32b Fully Flexible Bit-Precision and Vector Lockstep Execution Mode

A. Garofalo, G. Ottavi, A. Di Mauro, F. Conti, L. Benini, D. Rossi
DEI

University of Bologna, Italy
angelo.garofalo@unibo.it

September 6-9, 2021

- **Introduction & Motivation**
- Dustin Architecture Overview
 - Tunable Mixed-Precision Computation
 - Vector Lockstep Execution Mode
- Chip Results Summary
- Comparison with the State-of-the-art
- Conclusion

Introduction & Motivation

□ Extreme Edge AI and TinyML

- Low latency and network load compared to cloudML;
- Eases privacy concerns

□ Challenges

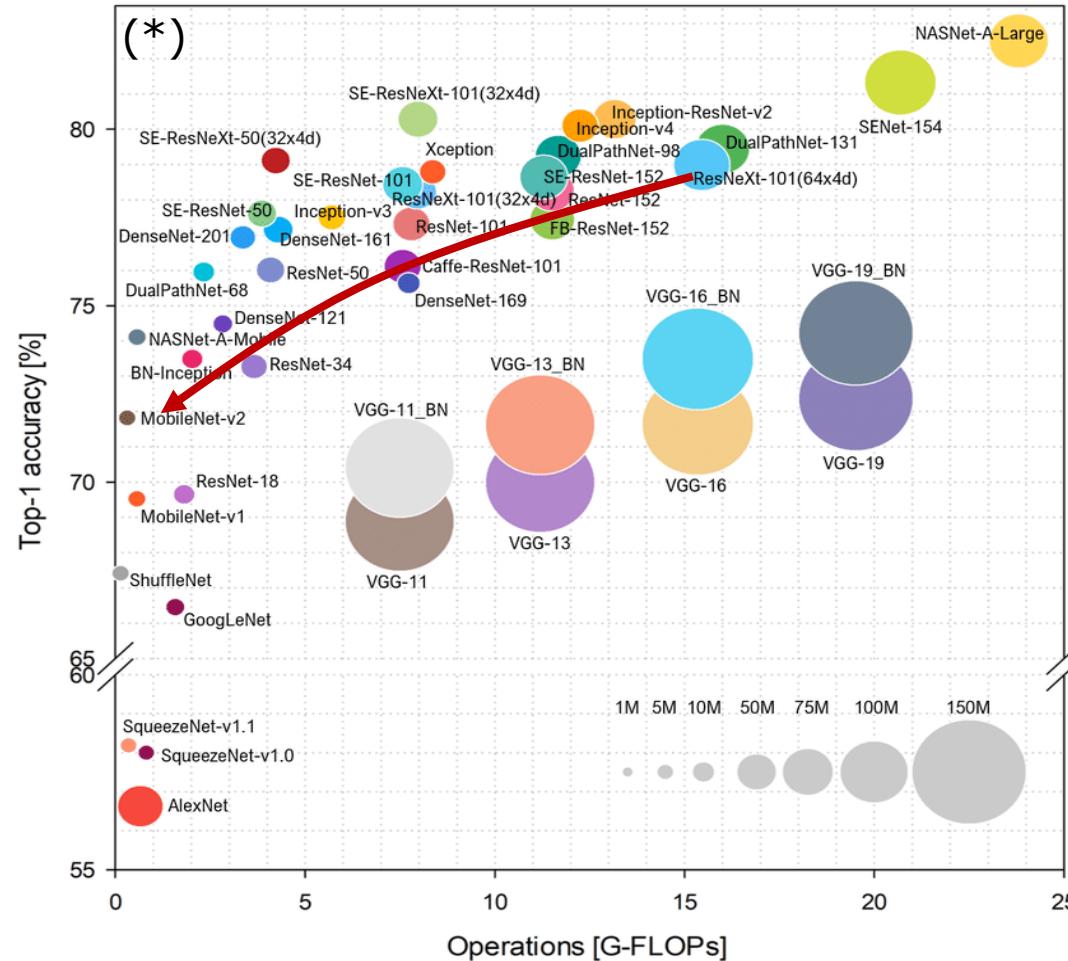
- High computational and memory requirements (ML + DL);
- Limited resources on IoT End-Nodes (microcontrollers).

□ Opportunities

- Reduce DL/ML model size;
- Low-Bitwidth Mixed-Precision computation.
- Reduce instruction fetch and decode overhead exploiting data-parallel execution



Quantized Neural Networks (QNNs)



Quantization Method	Top1 Accuracy	Weight Memory Footprint
Full-Precision	70.9%	16.27 MB
INT-8	70.1%	4.06 MB
INT-4	66.46%	2.35 MB
Mixed-Precision	68%	2.09 MB

Mixed-precision Quantized Neural Networks (QNNs) are the natural target for execution on constrained edge platforms.

(*) Bianco, Simone, Remi Cadene, Luigi Celona, and Paolo Napoletano. "Benchmark analysis of representative deep neural network architectures." IEEE Access 6 (2018): 64270-64277.

Edge AI Computing Platforms

	ASICs	FPGAs	MCUs
Throughput [Gop/s]	1 K – 50 K	10 – 200	0.1 – 2
Energy Efficiency [Gop/s/W]	10 K – 100 K	1 - 10	1 – 50
Flexibility	Low	Medium	High
Cost	High	Medium	Low

☐ IoT End-Nodes scenario:

- Must be inexpensive and software programmable (MCUs);
- SoA ISAs (RISC-V (*), ARM (**)) support only integer uniform arithmetic (with SIMD);
- Huge overhead to perform mixed-precision computation for data casting and packing;

☐ This Work:

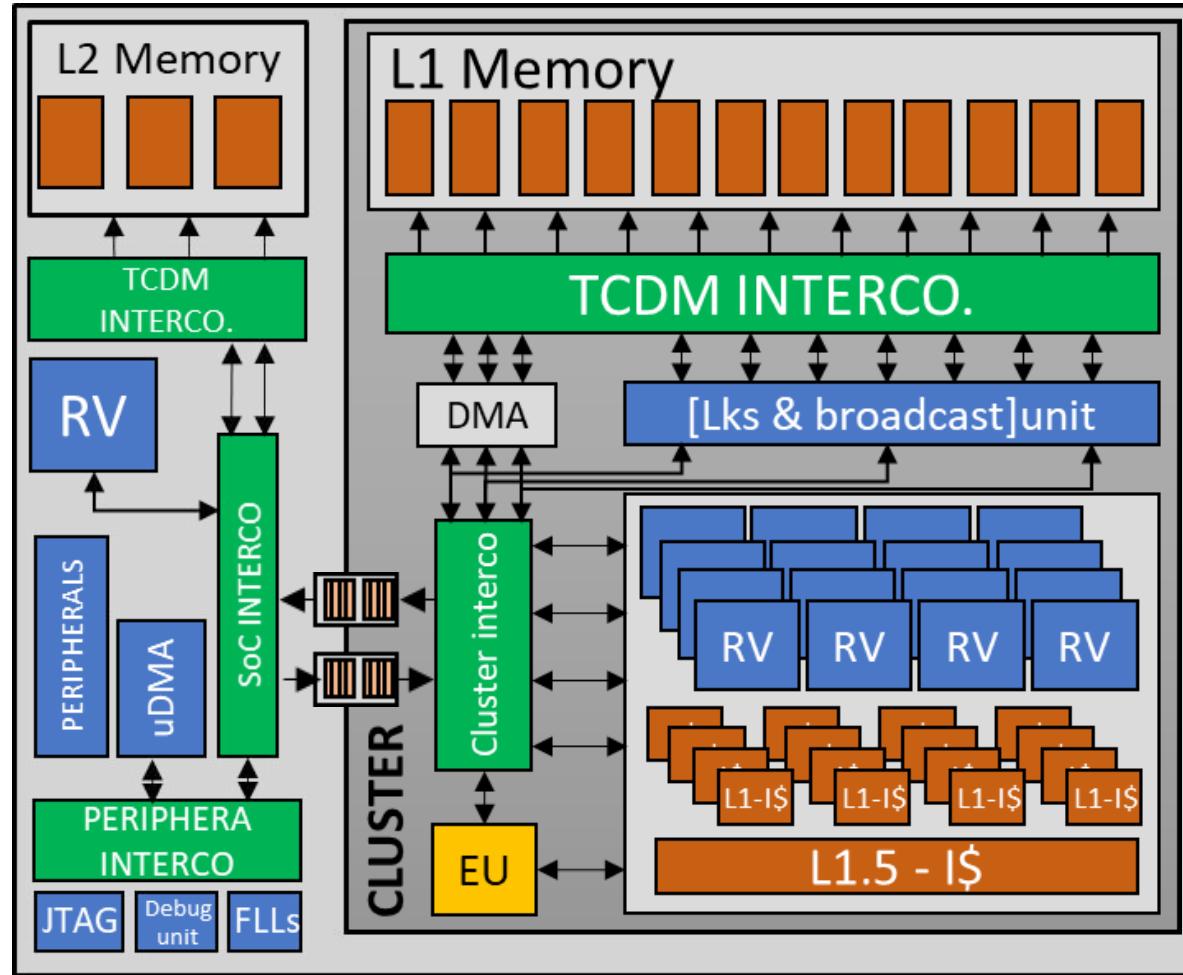
- Low-Power IoT End-Node with a fully programmable RISC-V accelerator cluster;
- Mixed precision 2b-to-32b SIMD instructions in the ISA of the cores;
- *Vector Lockstep Exec. Mode* to boost the Efficiency on data-parallel DL/ML algorithms.

(*) Garofalo et al. "XpulpNN: Enabling Energy Efficient and Flexible Inference of Quantized Neural Networks on RISC-V based IoT End Nodes." IEEE Transactions on Emerging Topics in Computing (2021).

(**) D. E. Joseph Yiu, "Introduction to the arm cortex-m55 processor. : <https://pages.arm.com/cortex-m55-introduction.html>," Feb. 2020.

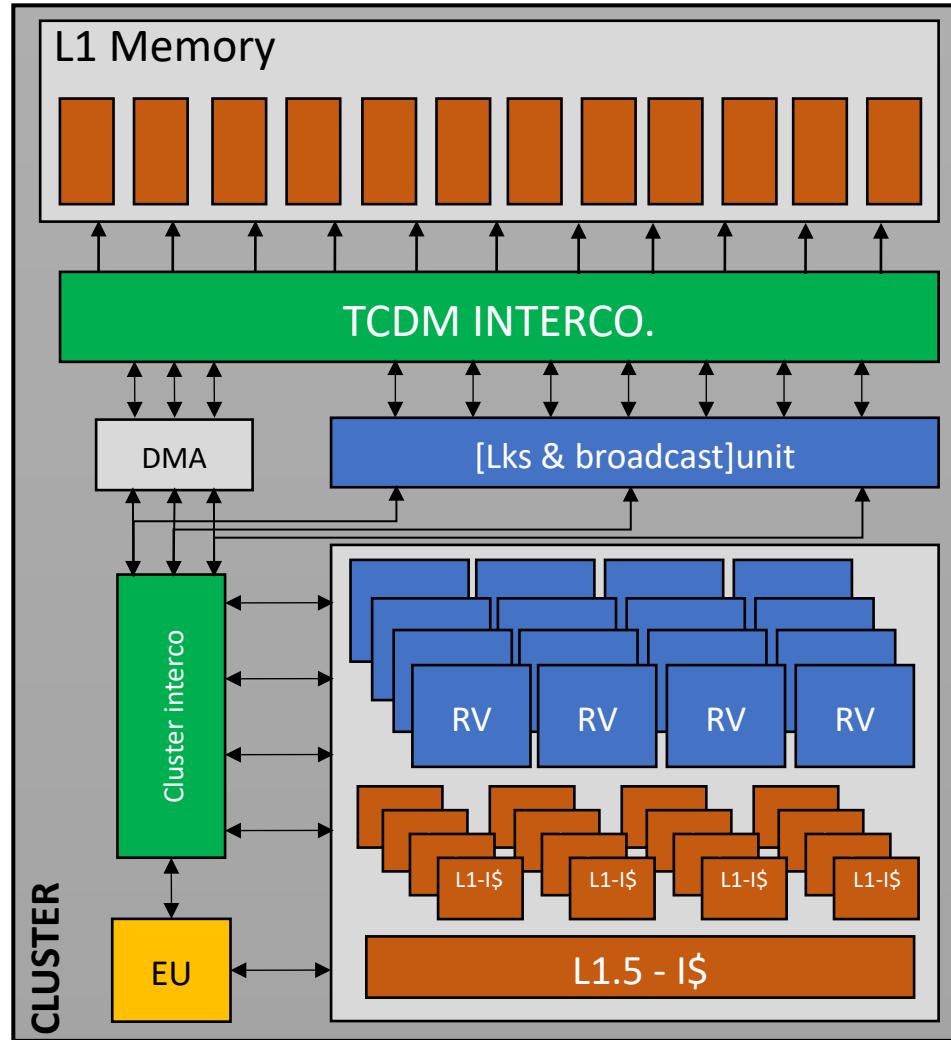
- Introduction & Motivation
- **Dustin Architecture Overview**
 - Tunable Mixed-Precision Computation
 - Vector Lockstep Execution Mode
- Chip Results Summary
- Comparison with the State-of-the-art
- Conclusion

Dustin: Architecture Overview



- **The SoC is a Low-Power IoT End-Node with AI edge computing capabilities**
 - Microcontroller
 - 1 RISCV core
 - 112 kB of L2 memory
 - Rich sets of peripherals (UART, I2C, CAM ift..)
 - JTAG (Debug), GPIOs, ROM
 - Interrupt Controller
 - Parallel cluster accelerator of fully programmable RISC-V cores

Dustin: Cluster

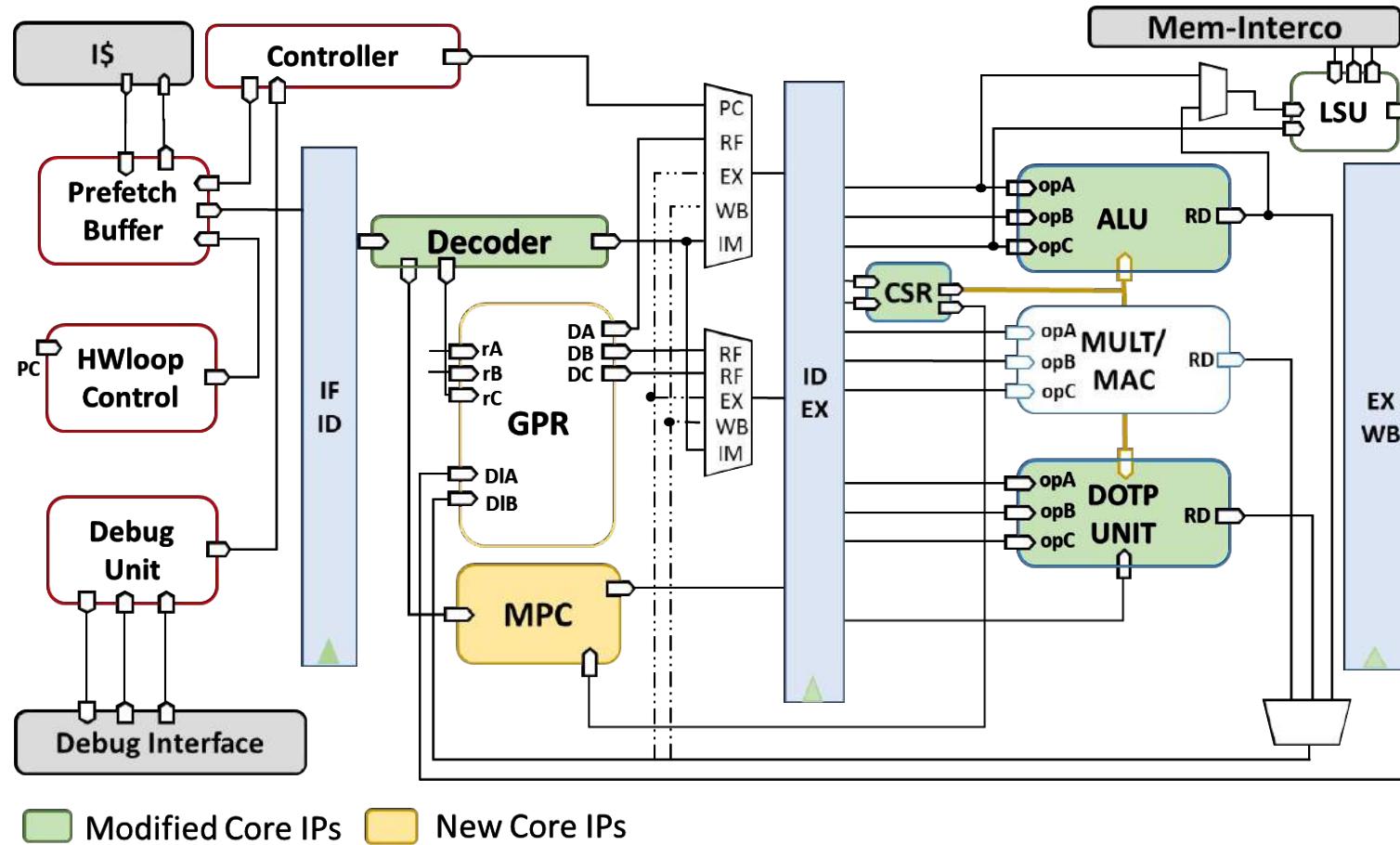


□ Accelerator Cluster

- 16 RI5CY (*) cores augmented with 2b-to-32b SIMD instructions;
- Software Configurable Vector Lockstep Execution Mode (VLEM);
- Single-cycle latency TCDM interco. leveraging a req/gnt protocol, word-level interleaved scheme.
- 128 kB of Shared Tightly-Coupled L1 Data Memory;
- Hierarchical Instruction Cache;
- High performance DMA (L2 <-> L1);
- Event Unit supporting efficient synchronization among the cores;

(*) Gautschi et al., Near-threshold RISC-V core with DSP extensions for scalable IoT endpoint devices, IEEE VLSI, 2017.

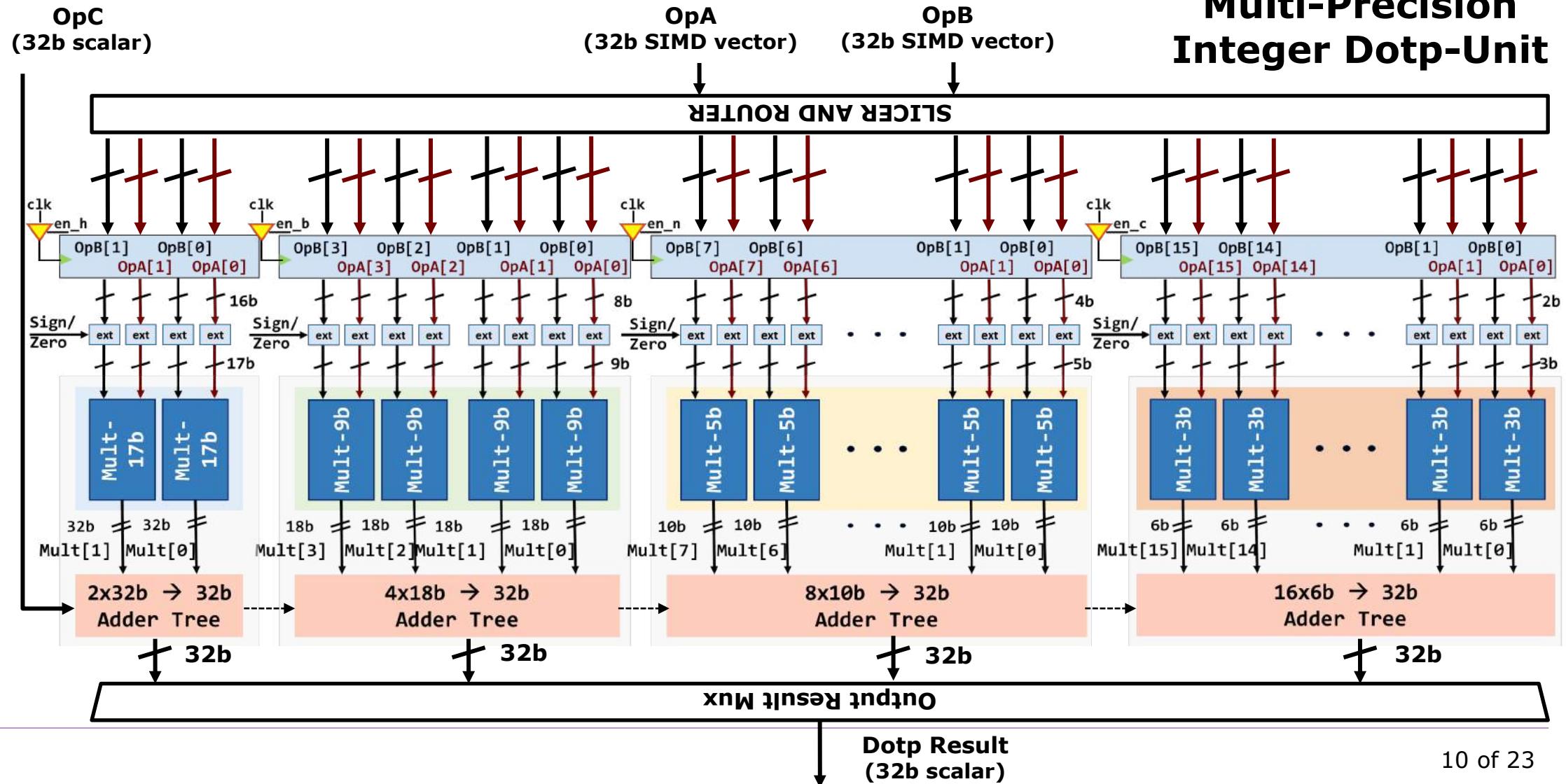
Core Enhancements



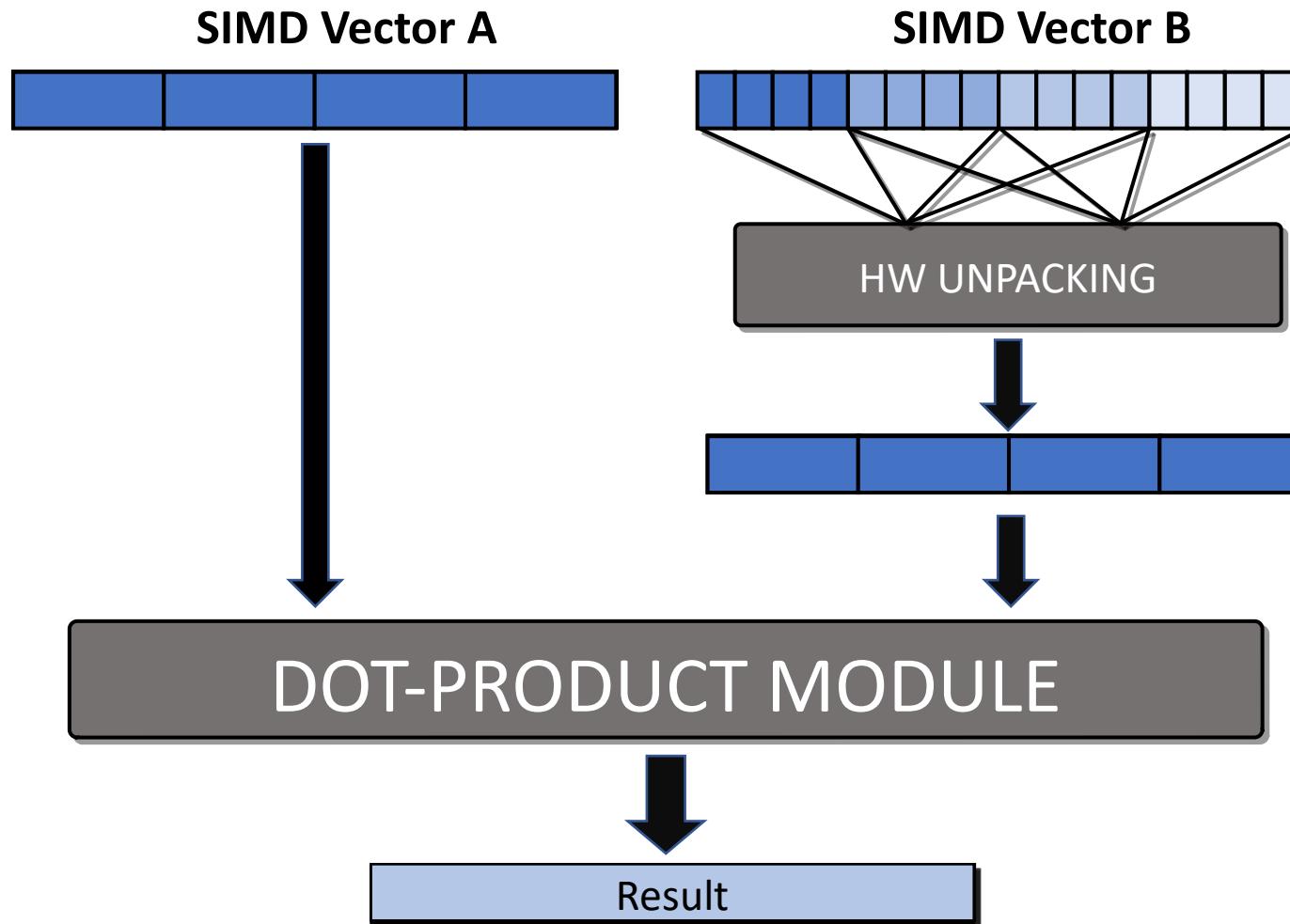
- **RI5CY:** 4-stage in order single-issue pipeline
- **ISA:** RV32IMCXPulpV2
- **XpulpV2 extensions:**
 - HW Loops;
 - Post-Increment LD/ST;
 - 16-/8-bit SIMD insns;
 - Bit Manip. insns.
- **Goal**
 - HW support for **mixed-precision** SIMD instructions;
- **Challenge**
 - Enormous number of instructions to be encoded in the ISA;
- **Solution**
 - **Dynamic Bit-Scalable Precision**

Extended Dot-Product Unit

Multi-Precision Integer Dotp-Unit



Mixed-Precision Controller



- Mixed-Precision operations require a controller: selection of the correct sub-word of the lowest precision operand (Vector B) to be used in current SIMD op.
- The controller is programmed by control status registers (CSRs).

Dynamic Bit-Scalable Execution

Standard Instructions						Virtual Instructions	
<code>pv.dotsp.h</code>						<code>pv.dotsp.v</code>	<code>pv.sdotsp.v</code>
<code>pv.dotsp.b</code>						<code>pv.dotsp.sc.v</code>	<code>pv.sdotsp.sc.v</code>
<code>pv.dotsp.n</code>						<code>pv.dotsp.sci.v</code>	<code>pv.sdotsp.sci.v</code>
<code>pv.dotsp.c</code>						<code>pv.dotup.v</code>	<code>pv.sdotup.v</code>
<code>pv.dotsp.m4x2</code>						<code>pv.dotup.sc.v</code>	<code>pv.sdotup.sc.v</code>
<code>pv.dotsp.m8x2</code>						<code>pv.dotup.sci.v</code>	<code>pv.sdotup.sci.v</code>
<code>pv.dotsp.m8x4</code>						<code>pv.dotusp.v</code>	<code>pv.sdotusp.v</code>
<code>pv.dotsp.m16x8</code>						<code>pv.dotusp.sc.v</code>	<code>pv.sdotusp.sc.v</code>
<code>pv.dotsp.m16x4</code>						<code>pv.dotusp.sci.v</code>	<code>pv.sdotusp.sci.v</code>
<code>pv.dotsp.m16x2</code>							
<code>pv.dotsp.sc.h</code>							
<code>pv.dotsp.sc.b</code>							
<code>pv.dotsp.sc.c</code>							
<code>pv.dotsp.sc.n</code>							
<code>pv.dotsp.sc.m4x2</code>							
<code>pv.dotsp.sc.m8x2</code>							
<code>pv.dotsp.sc.m8x4</code>							
<code>pv.dotsp.sc.m16x8</code>							
<code>pv.dotsp.sc.m16x4</code>							
<code>pv.dotsp.sc.m16x2</code>							
<code>pv.dotsp.sci.h</code>							
<code>pv.dotsp.sci.b</code>							
<code>pv.dotsp.sci.c</code>							
<code>pv.dotsp.sci.n</code>							
<code>pv.dotsp.sci.m4x2</code>							
<code>pv.dotsp.sci.m8x2</code>							
<code>pv.dotsp.sci.m8x4</code>							
<code>pv.dotsp.sci.m16x8</code>							
<code>pv.dotsp.sci.m16x4</code>							
<code>pv.dotsp.sci.m16x2</code>							
<code>pv.packlo.b x15, x5, x6</code>							
<code>pv.packhi.b x15, x7, x8</code>							
<code>pv.sdotsp.b x20, x15, x10</code>							

instruction per format and type

Virtual Instructions

- No execution overhead
- Reuse of standard instruction sets

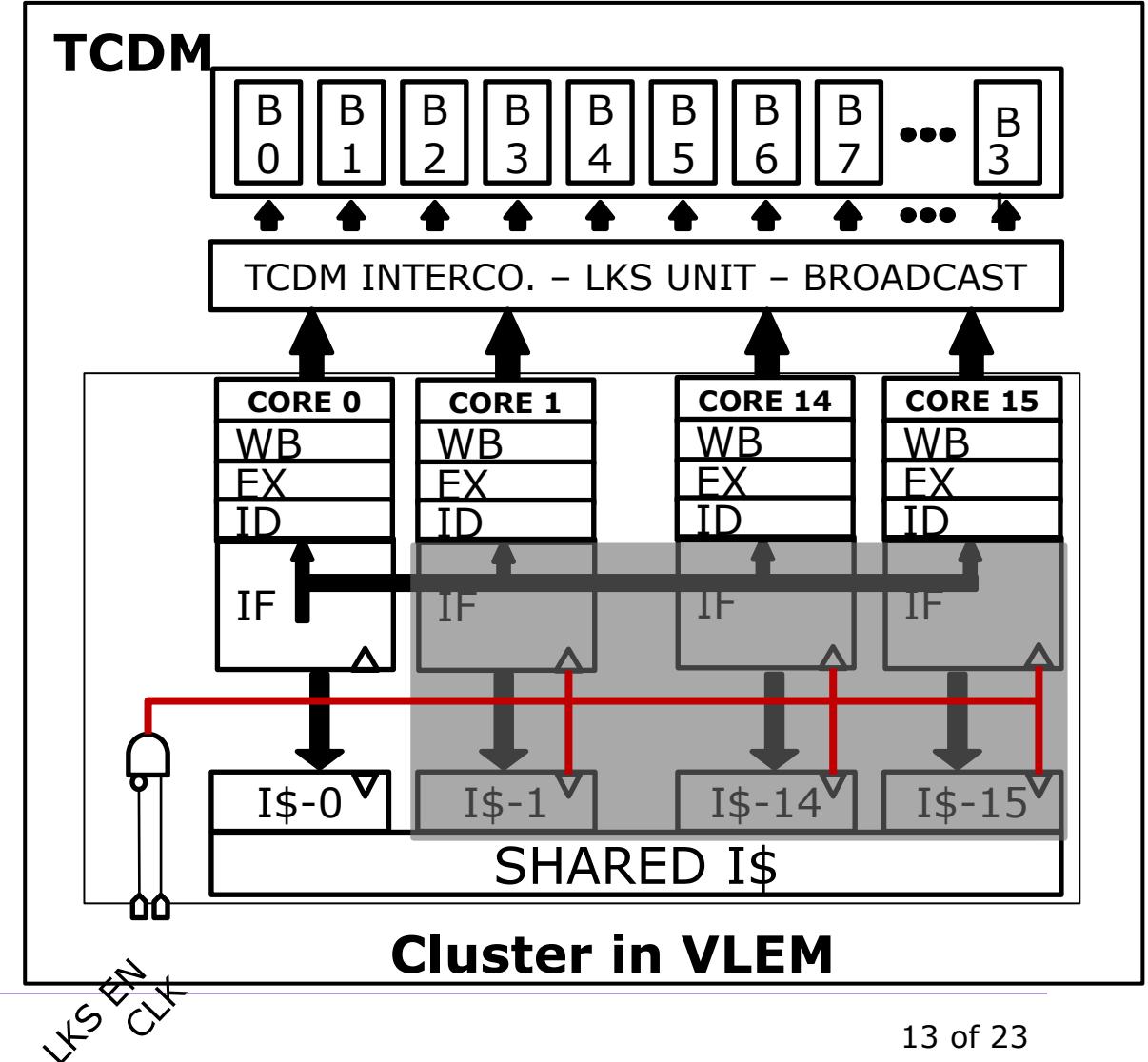
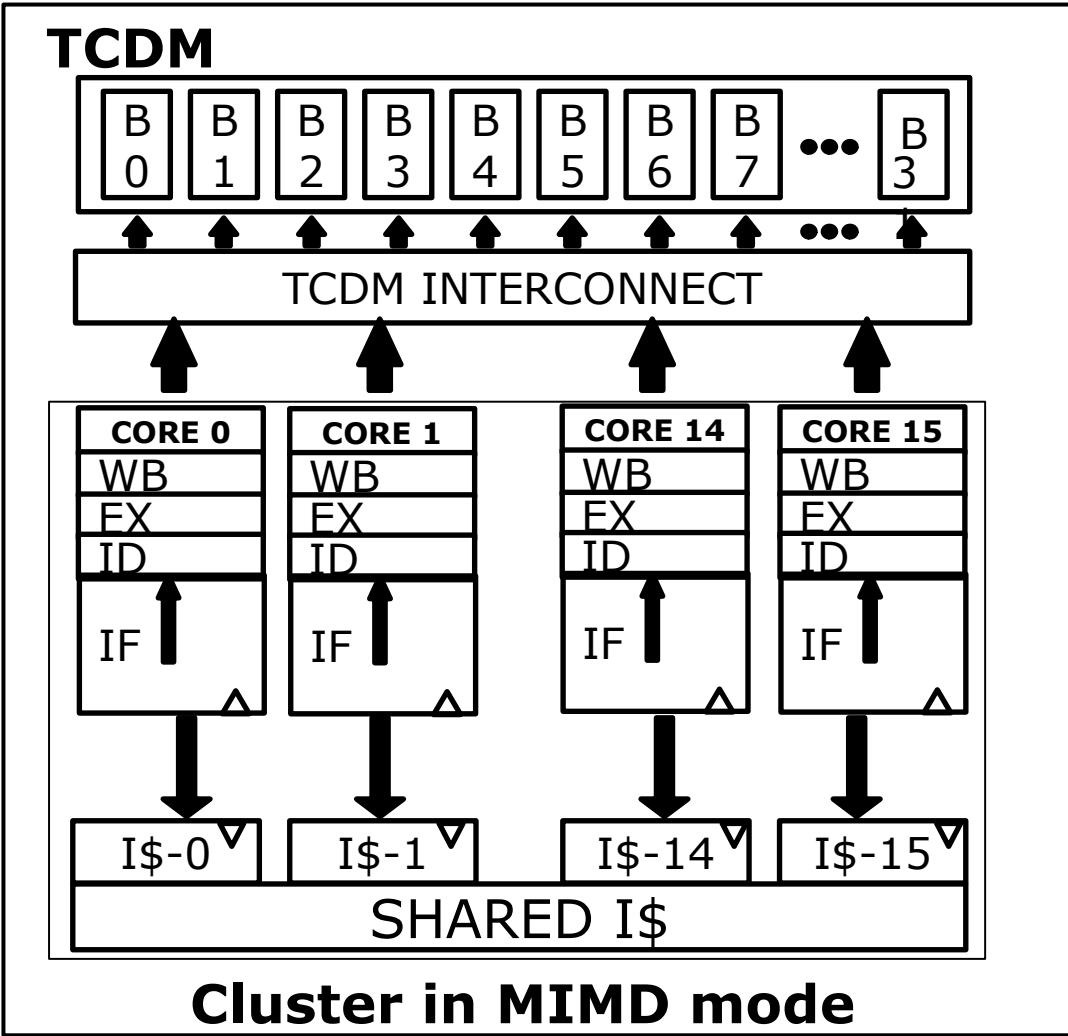
Dynamic Bit-Scalable Execution

10,4(x4!)
11,4(x5!)
p.v x20,x11,x10

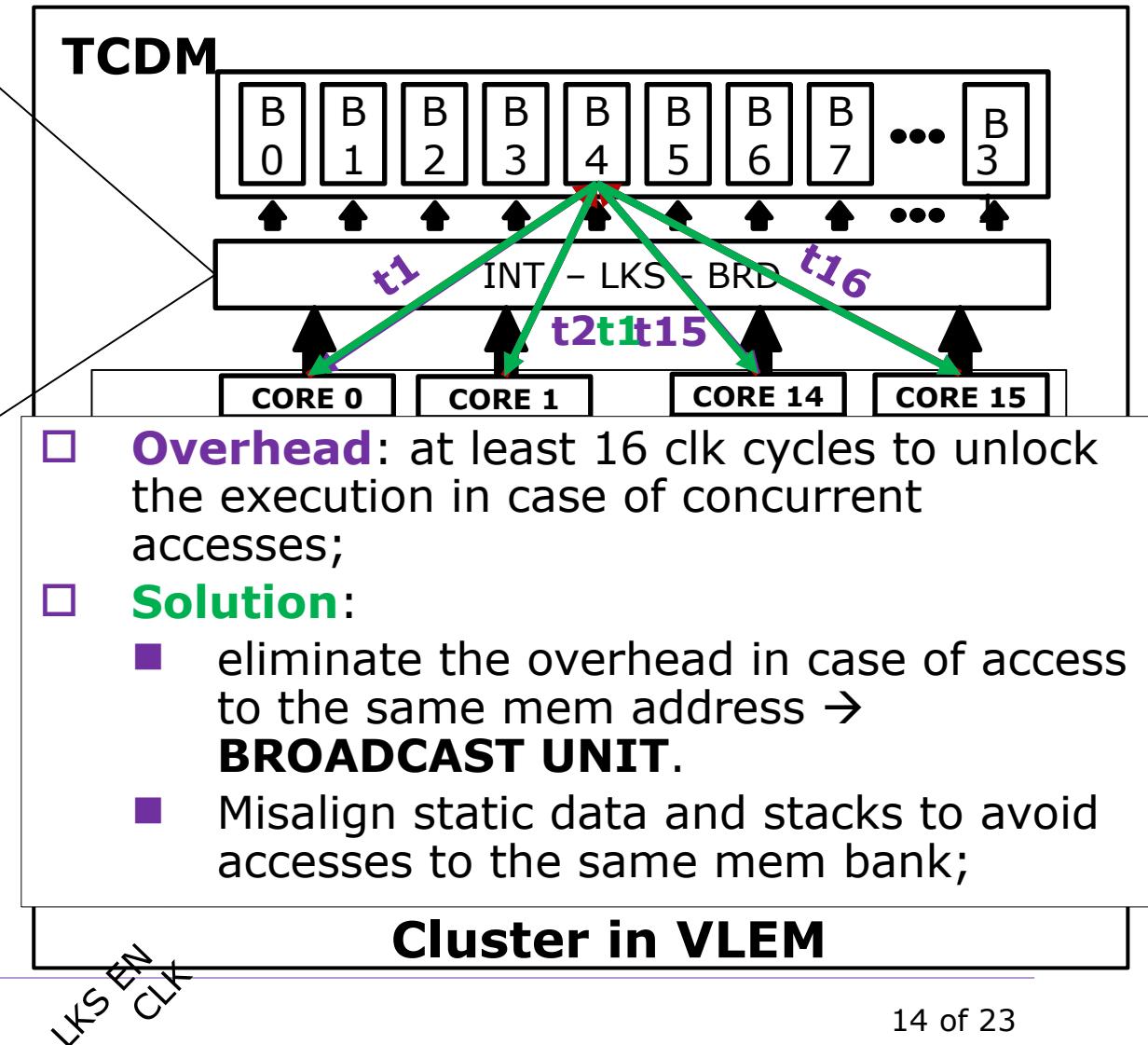
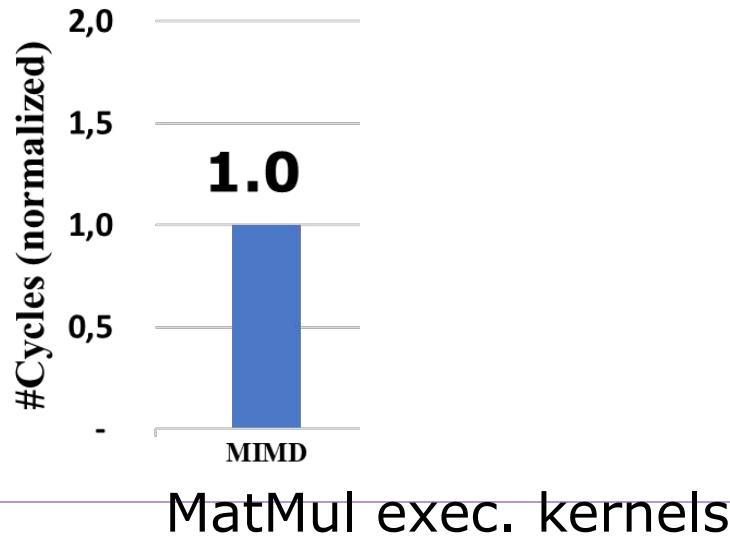
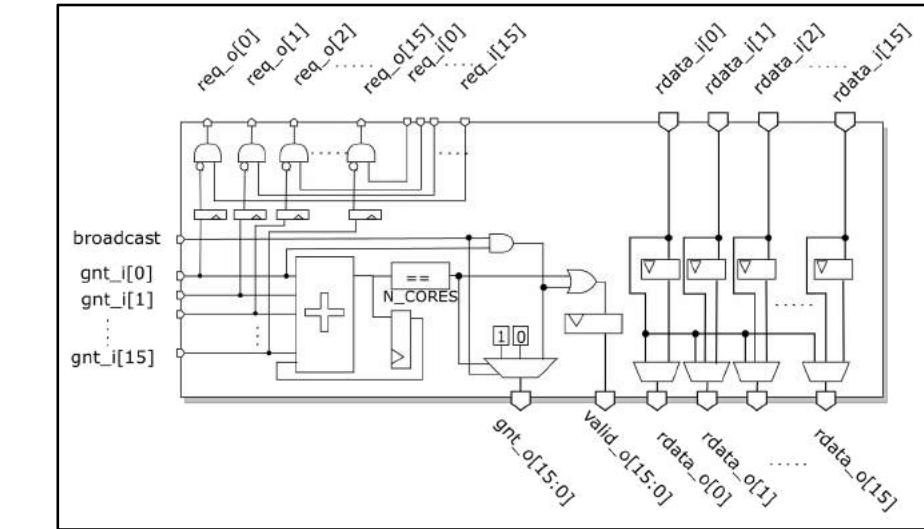
one instruction per type

```
int main()
{
    ....
    SIMD_FMT(M8x4);
    convolution(A, W, Res);
    ....
    SIMD_FMT(M8x2);
    convolution(A, W, Res);
    ....}
```

Vector Lockstep Exec. Mode (VLEM)



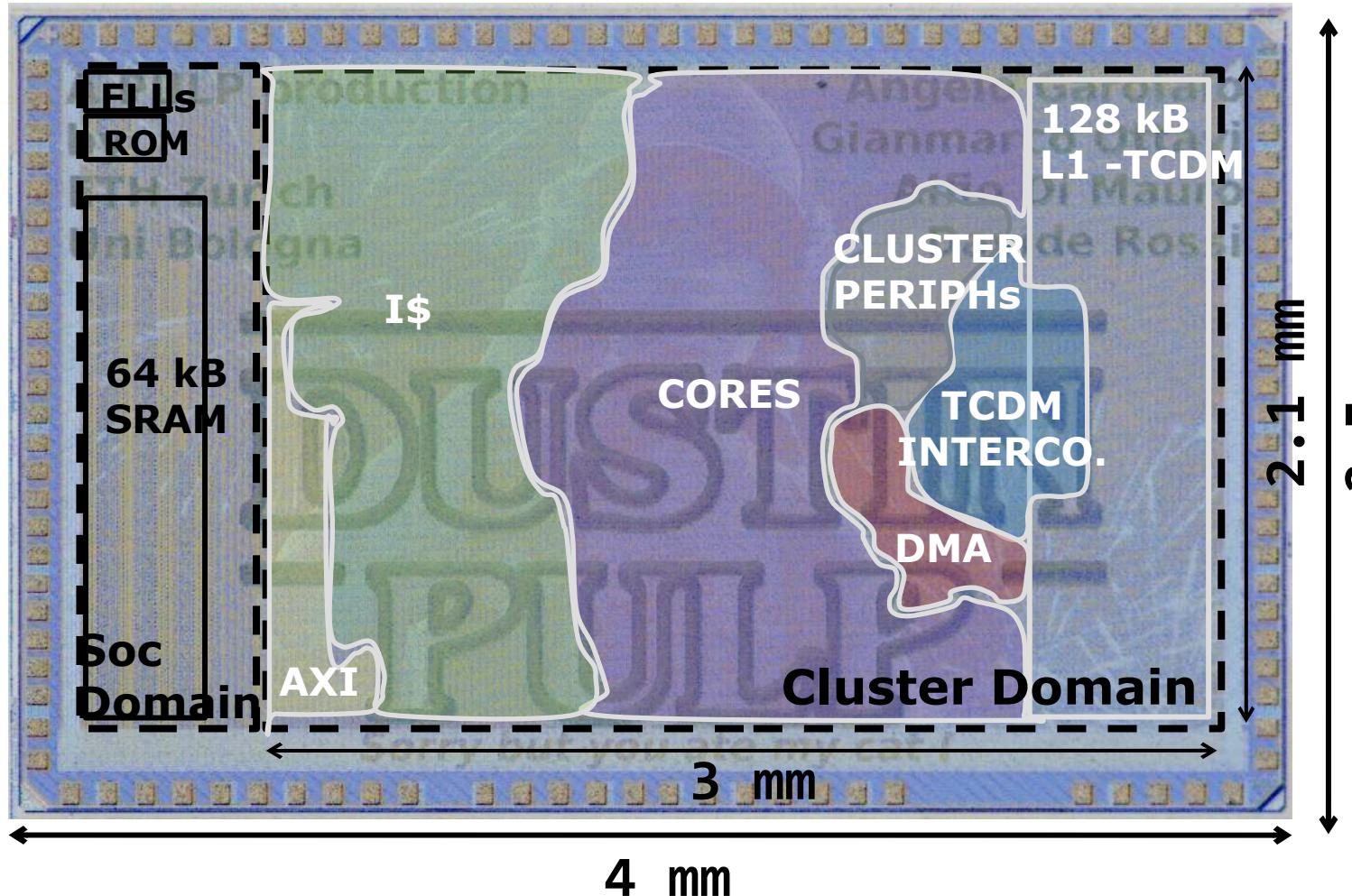
VLEM: Broadcast Unit



Outline

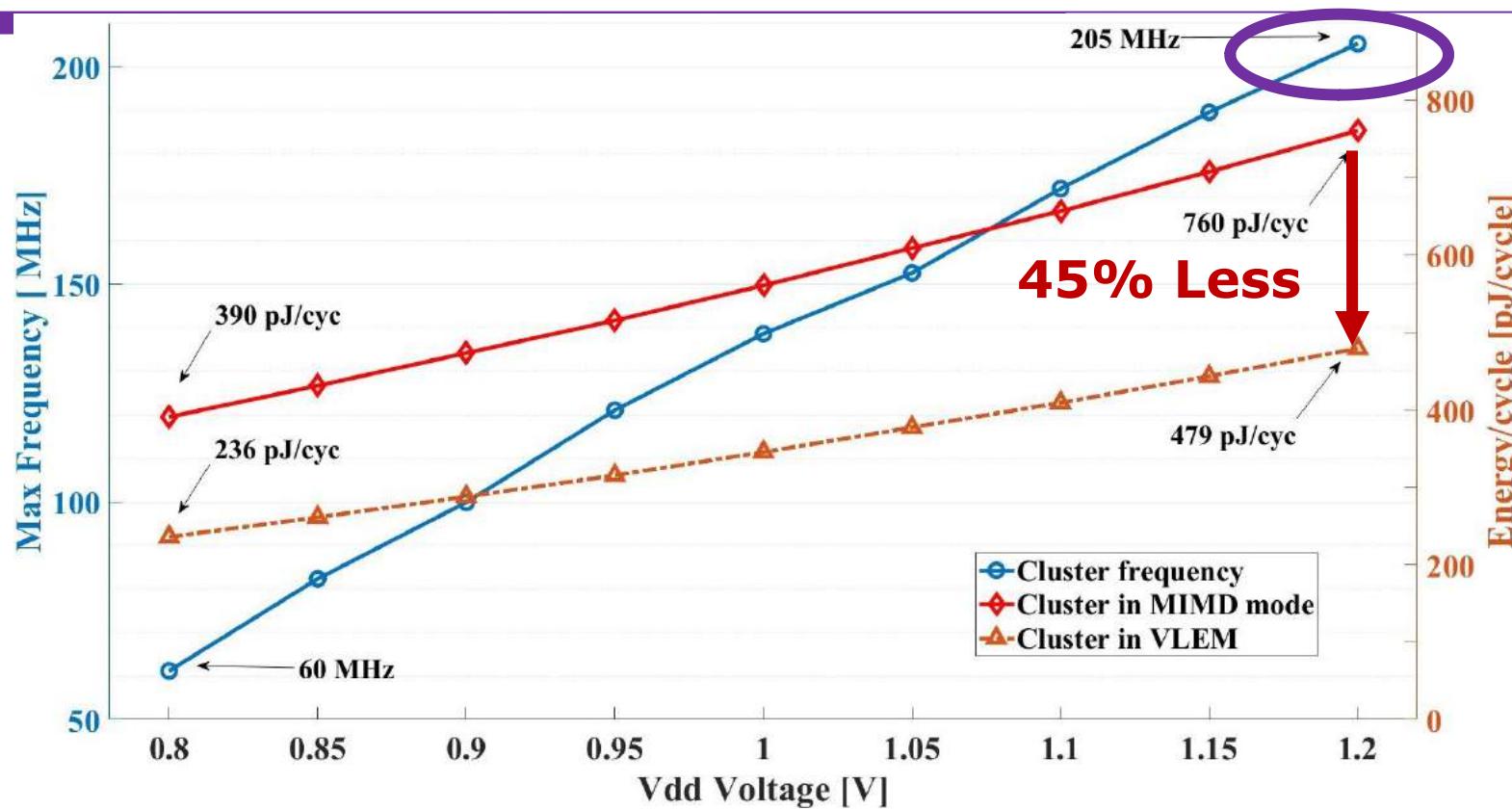
- Introduction & Motivation
- Dustin Architecture Overview
 - Tunable Mixed-Precision Computation
 - Vector Lockstep Execution Mode
- **Chip Results Summary**
- Comparison with the State-of-the-art
- Conclusion

Chip Results Summary



Technology	CMOS TSMC 65nm
Chip area	10 mm ²
Total SRAM	208 kB
VDD range	0.8V - 1.2V
Power Envelope	156 mW
Frequency Range	60 – 205 MHz
Integer Performance (8-bit)	15 GOPS
Integer Efficiency (8-bit)	303 GOPS/W
Integer Performance (4-bit)	30 GOPS
Integer Efficiency (4-bit)	562 GOPS/W
Peak Performance	58 GOPS
Peak En. Efficiency	1.15 TOPS/W

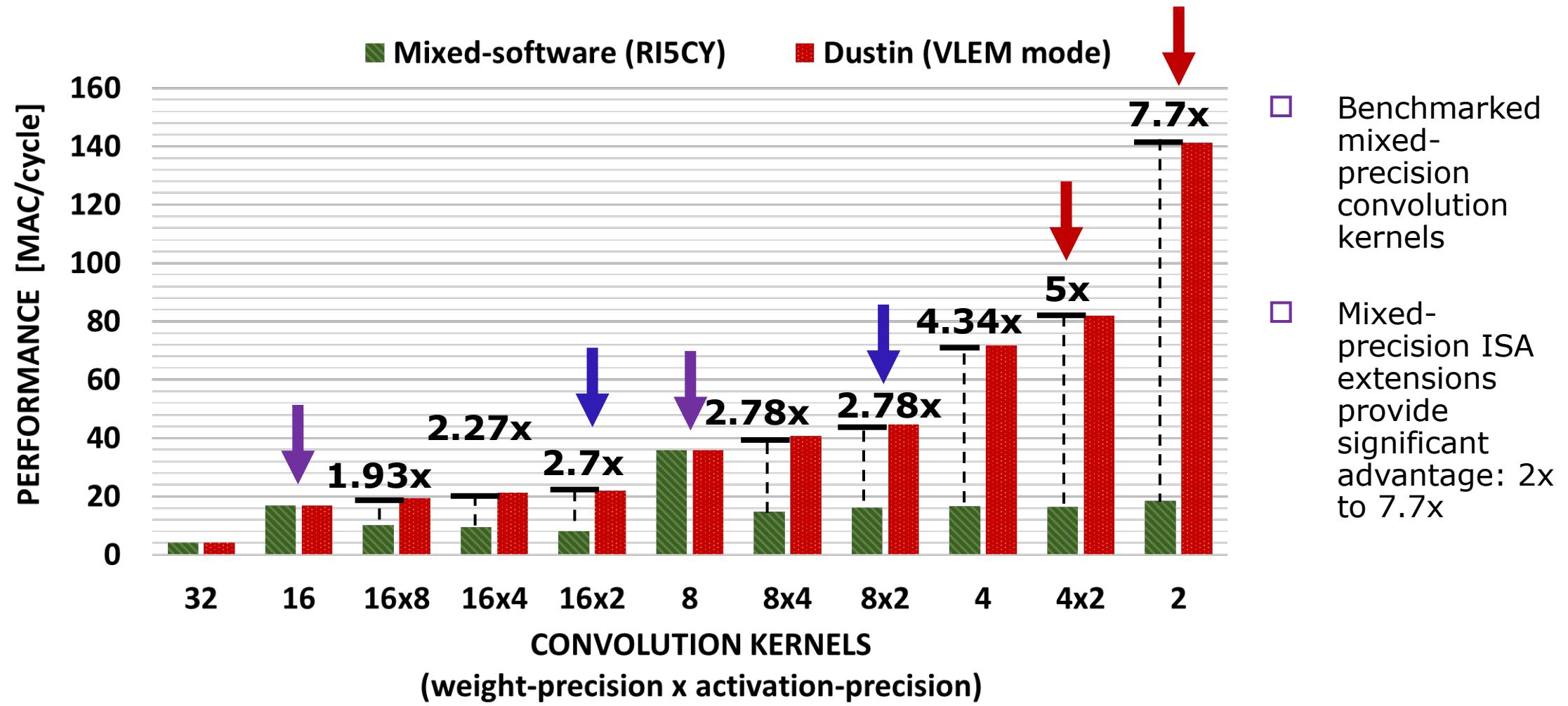
Voltage vs. Frequency



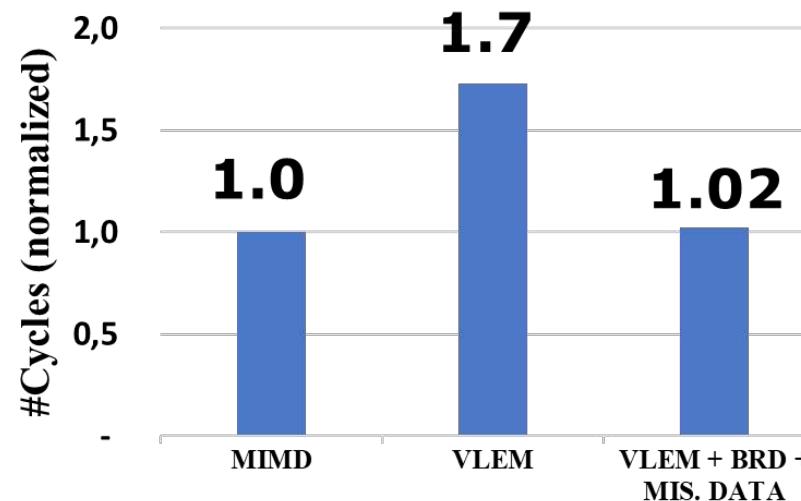
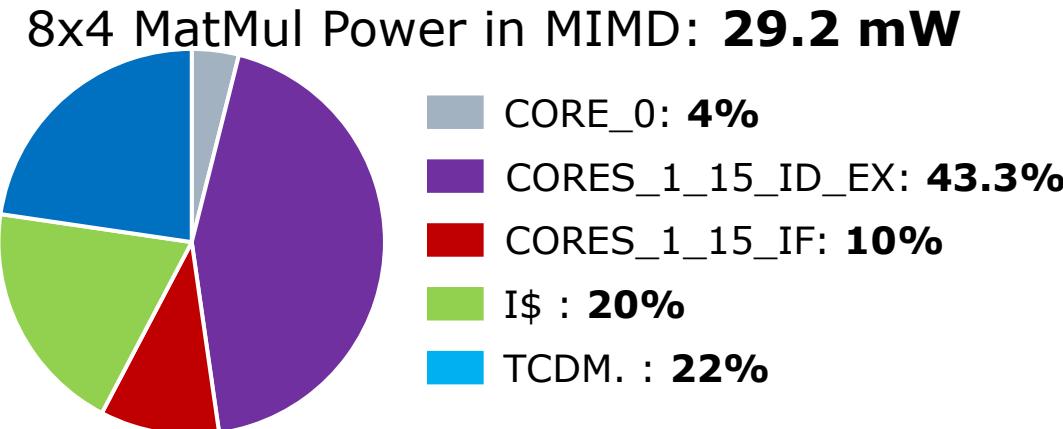
**VLEM, no impact
on max
operating freq.
of the cluster**

- Measurements of the Cluster;
- Maximum frequency 205 MHz @ 1.2 V;
- ~45% energy saving** with the Cluster in VLEM wrt MIMD mode.

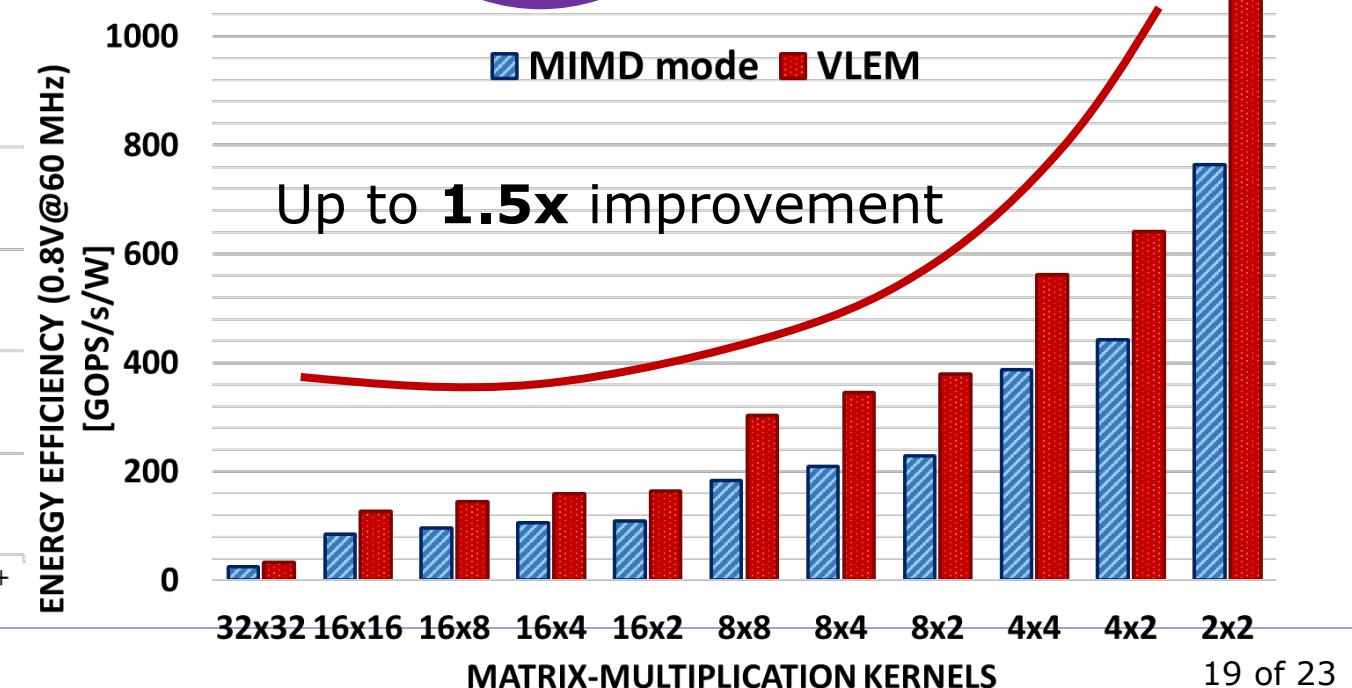
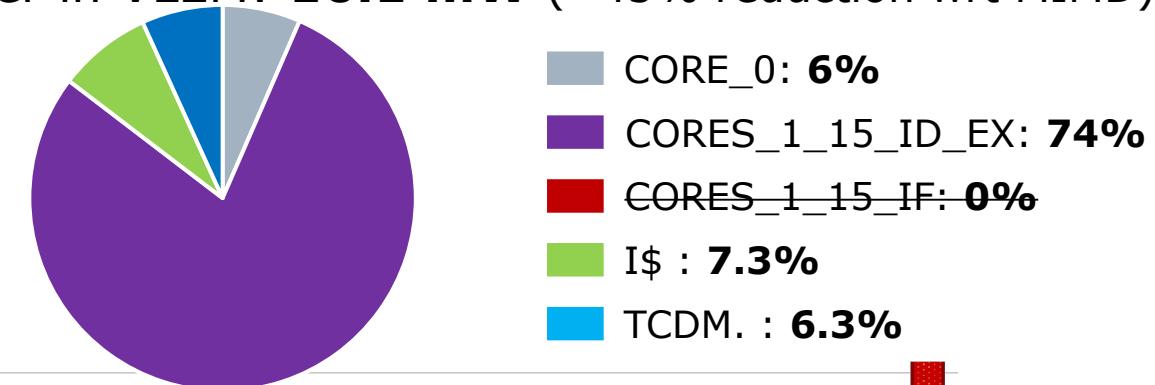
Performance on MatMul kernels



Energy Efficiency on MatMul kernels



Power in VLEM: **16.1 mW** (~45% reduction wrt MIMD)



Outline

- Introduction & Motivation
- Dustin Architecture Overview
 - Tunable Mixed-Precision Computation
 - Vector Lockstep Execution Mode
- Chip Results Summary
- **Comparison with the State-of-the-art**
- Conclusion

Comparison with the SoA

	SleepRunner [6]	SamurAI [7]	Mr.Wolf [8]	Vega [9]	Dustin (this work)
Technology	CMOS 28nm FDSOI	CMOS 28nm FDSOI	CMOS 40nm LP	CMOS 22nm FDSOI	CMOS 65nm
Die Area	0.68 mm ²	4.5 mm ²	10 mm ²	12 mm ²	10 mm ²
Applications	IoT GP	IoT GP + DNN	IoT GP + DNN	IoT GP + NSA+DNN	IoT GP + DNN + QNNs
CPU/ISA	CM0DS Thumb-2 subset	1x RI5CY RVC32IMFXpulp	9 x RI5CY RVC32IMFXpulp	10 x RI5CY RVC32IMFXpulp+SF	16 x MPIC CORES (RISC-V)
Int Precision (bits)	32	8, 16, 32	8, 16, 32	8, 16, 32	2, 4, 8, 16, 32 (plus Mixed-Precision)
Supply Voltage	0.4 - 0.8 V	0.45 - 0.9 V	0.8 - 1.1 V	0.5 – 0.8 V	0.8 - 1.2 V
Max Frequency	80 MHz	350 MHz	450 MHz	450 MHz	205 MHz
Power Envelope	320 µW	96 mW	153 mW	49.4 mW	156 mW
¹Best Integer Performance	31 MOPS (32b)	1.5 GOPS (8b) ²	12.1 GOPS (8b)	15.6 GOPS (8b)	15 GOPS (8b) 30 GOPS (4b) 58 GOPS (2b)
¹Best Integer Efficiency	97 MOPS/mW @ 18.6 MOPS (32b)	230 GOPS/W @110 MOPS (8b) ²	190 GOPS/W @ 3.8 GOPS (8b)	614 GOPS/W @ 7.6 GOPS	303 GOPS/W @4.4 GOPS (8b) 570 GOPS/W @8.8 GOPS (4b) 1152 GOPS/W @17.3 GOPS(2b)

¹ 2 OPs = 1 8-bit (or 4-bit or 2-bit) MAC on MatMul benchmark unless differently specified.

² For fair comparison we consider the execution on software programmable cores.

Dustin supports Mixed-Precision computation in HW

Better efficiency wrt solutions in 28nm and 40nm tech node

Comparable efficiency wrt Vega (22 nm)

Outline

- Introduction & Motivation
- Dustin Architecture Overview
 - Tunable Mixed-Precision Computation
 - Vector Lockstep Execution Mode
- Chip Results Summary
- Comparison with the State-of-the-art
- **Conclusion**

Conclusion

- Dustin SoC: IoT end-node with AI computing capabilities in **tsmc 65 nm** tech node;
- RISC-V cores featuring 2b-to-32b bit-precision instruction set architecture (ISA) extensions enabling fine-grain tunable mixed-precision computation (**2x** to **7x** speed-up w.r.t. RI5CY);
- Software reconfigurable cluster in Vector Lockstep Execution Mode (**~45% energy saving** w.r.t. MIMD mode);
- Dustin is competitive with IoT end-nodes using much more scaled technology nodes (Peak Perf. **58 GOPS**, Peak Eff. **1.15 TOPS/W**).
- Despite a less scaled tech node, we reach energy efficiency in the order of **TOPS/W** → Comparable with ASIC solutions.