#### Is an AI Accelerator All You Need? Overcoming Amdahl's Law With Tightly-Coupled Specialized Accelerators

Angelo Garofalo ETH Zürich, Switzerland & University of Bologna, Italy

ISSCC 2023, FORUM 6 "The Future of Heterogeneous Multi-Core Architectures for Al and Other Specialized Processing"

## Outline

- □ The quest for energy-efficient AI at the Edge of the IoT
  - The performance-power proportionality goal
  - Dealing with the end of the Dennard's scaling
- □ Low-power heterogeneous computing: mitigating deep acceleration effects
  - Multi-core heterogeneous systems
  - Tightly-coupled digital acceleration
  - Heterogeneous analog in-memory computing (AIMC) systems
- □ Insights for scaled-up heterogeneous powerful architectures
  - AIMC many-core architecture
  - Asymmetric chiplet-based systems
- Conclusions & Outlook

Angelo Garofalo

## The Quest for Advanced AI Computing at the Edge



Computing requirements for Edge AI

- AI capabilities under tight power (passive cooling, few <u>hundreds mW</u> of budget)
- Real-time constraints (<u>Latency</u> as important as <u>Throughput</u>)
- Tight area budget (on-chip memory <u>limited</u> to ~1MB)

Angelo Garofalo

## **Reduced Precision DL Models**



#### Quantized MobileNet-v1 model

Quantization Method	Top1 Accu	racy	Weight Memory Footprint				
Full-Precision	70.9%		16.27 MB				
INT-8	70.1%	0.8%	4.06 MB	<b>4</b> x			
INT-4	66.46%	4.4%	2.35 MB	<b>7</b> x			
Mixed-Precision	68%	<b>*</b> 2.9%	2.09 MB	<b>▼ 8x</b>			

[ Rusci et al., PML&S (2020)]

- □ MVM workload dominates (massive parallelism)
- □ Strong resiliency to Quantization (<u>low-precision</u> integer arithmetic is ok)
- □ Mixed-precision quantized neural networks (QNNs) are the natural target for Edge AI

Angelo Garofalo

#### The Performance-Power Proportionality Goal



## Technology Scaling is Not Enough



Increase number of processing cores (..but up to how many cores?)
 Hardware specialization (..but what about flexibility?)

Angelo Garofalo

## HP Systems: Tesla Full Self Driving SoC Example



Full Self-Driving (FSD) SoC in 14nm FINFET technology
 6 B Transistors within 265mm<sup>2</sup> die area, ~100W Power Envelope
 68 GB/s on-chip BW

Angelo Garofalo

## Outline

- □ The quest for energy-efficient AI at the Edge of the IoT
  - The performance-power proportionality goal
  - Dealing with the end of the Dennard's scaling
- □ Low-power heterogeneous computing: mitigating deep acceleration effects
  - Multi-core heterogeneous systems
  - Tightly-coupled digital acceleration
  - Heterogeneous analog in-memory computing (AIMC) systems
- □ Insights for scaled-up heterogeneous powerful architectures
  - AIMC many-core architecture
  - Asymmetric chiplet-based systems
- Conclusions & Outlook

Angelo Garofalo

## Near-Threshold Parallel Computing



- □ SIMD insns rather than super-scalar execution for high energy-efficiency
- More PEs at the optimum energy point working on parallel workloads
- □ .. To achieve high performance-power proportionality (1pJ/Op @ GOPS)

Angelo Garofalo

## Parallel Ultra-Low Power (PULP) Platform



[Rossi et al. JSSC 2022]

- □ *N* Specialized RISC-V cores (domain-specific ISA instructions)
- Single-cycle latency Logarithmic Interconnect towards TCDM
- SW-managed (DMAs) tightly-coupled data mem (TCDM) rather than data caches
- Synchronization in HW for low-energy thread dispatching and PEs synch

Angelo Garofalo

## **Processor Specialization**



- □ [Baseline] RV32IMC: not good for machine learning workload (~44kGE)
- **[DSP] Xpulp**: 16/8-bit SIMD, HW loops, post-modified LD/ST & bit-manip. insns (~55kGE)
- **[AI] XpulpNN**: 4-bit, 2-bit SIMD (and mixed-precisions) & fused MAC-Load ops (~65kGE)

Angelo Garofalo

## **DNN Acceleration Through Domain-Specific ISA**



□ Power ① (slightly); Latency  $\mathbb{Q} \oplus \mathbb{Q}(>10x)$  on QNN kernels  $\Rightarrow$  Energy  $\oplus \oplus$ 

Parallelization of the workload (multi-core) for high throughput

Angelo Garofalo

## 8 RVXpulpnn Cores Cluster (22nm)



Angelo Garofalo

## Multi-Core Heterogeneous Computing



- Reduce unnecessary power consumption (not spent in computation)
- Exploit convolution's instruction and memory data access pattern regularity
- □ Increase energy efficiency at low extra-area cost
  - $\Box \rightarrow$  reconfigurable MIMD/SIMD architecture

Angelo Garofalo

## Multi-Core Heterogeneous Computing - 2



- □ Cores enter in SIMD (VLEM) mode when executing regular kernels (in two clock cycles)
- □ In SIMD, instruction flow orchestrated only by the MAIN core  $\rightarrow$  Less energy
- □ Cores resume in MIMD mode on divergent branches (..or control tasks)  $\rightarrow$  Flexibility

Angelo Garofalo

ISSCC 2023 - Forum 6.4: Is an AI accelerator all you need? Overcoming Amdahl's law with tightly-coupled heterogeneous accelerators. 15 of 44

## Multi-Core Heterogeneous Computing - 3



- **Overhead**: many clk cycles to unlock execution in case of concurrent accesses
  - Eliminate overhead to access at same address → BROADCAST UNIT
  - Misalign static data and stacks to avoid accesses to the same mem bank

Angelo Garofalo

## System Scalability

Results from convolution kernels run on 65nm silicon prototype



Performance bottleneck at the core's boundaries (single 32-bit data port)

- □ Energy efficiency bounded by limited scalability of low-latency local interco
- Custom datapaths to improve throughput and efficiency...
  - □ ..but at which cost?

Angelo Garofalo

## Specialized AI Accelerators - Digital



Custom MAC units, data-paths, data precision (binary, ternary, 2-8 bits)

- Perf. [0.1, 1] TOPS, En. Eff. [1, 100] TOPS/W
- □ Hardware challenge moved at the boundaries, on data movements

Angelo Garofalo

ISSCC 2023 - Forum 6.4: Is an AI accelerator all you need? Overcoming Amdahl's law with tightly-coupled heterogeneous accelerators.

18 of 44

## Specialized AI Accelerators - Analog



- Weights are stored where the computation happens (nvAIMC arrays)
- □ Analog MAC more efficient than digital MAC (10x better perf. and en. eff.)
- Energy efficiency vs. accuracy
- Compute density vs. re-programmability

Angelo Garofalo

ISSCC 2023 - Forum 6.4: Is an AI accelerator all you need? Overcoming Amdahl's law with tightly-coupled heterogeneous accelerators. 19 of 44

## Edge AI Processors



Angelo Garofalo

## **Towards Heterogeneous Integration**

- Performance bottleneck keeps shifting toward non-accelerated kernels
- □ (as consequence of Amdahl's effect)
  - Solutions
  - Balance workload through many acceleration sub-systems
  - Residual computation performed in SW
- □ Hardware challenge becomes efficient data movements within the system
  - Solutions
  - High-bandwidth low-latency interconnect
  - Reduce data-exchange latency through tightly-coupled integration schemes

## **Tightly-Coupled Digital Acceleration**



- Acceleration with flexibility: zero-copy HW/SW cooperation
- Many tightly-coupled accelerators for Amdahl's effects mitigation
- HCI: Low-Latency (one clk cyc), High-Bandwidth (76B/cyc)

Angelo Garofalo

## Hardware Processing Engines (HWPEs)



#### □ HWPE efficiency:

- Dedicated control (no I-fetch) with shadow registers (overlapped config-exec)
- Specialized high-BW interconnect towards L1 buffer (on data interface)
- Specialized compute datapath & minimal internal storage (e.g. line buffers)

Angelo Garofalo

## Heterogeneous Cluster Interconnect



#### Processors with their narrow memory ports are arbitrated in the LIC

Angelo Garofalo

## Heterogeneous Cluster Interconnect - 2



HWPE exposes a unified high-bandwidth port (e.g. 256-bit) towards the shallow interconnect (simple "wide" port reordering block without self-contention)

Angelo Garofalo

## Heterogeneous Cluster Interconnect - 3



□ Final level dispatches accesses to single 32-bit memory banks and manages conflicts by means of rotating configurable priority (e.g. max 3 stall cycles)

□ If no conflicts, HWPEs and cores can access concurrently the memory.

Angelo Garofalo

## HW Synchronizer



- □ Fast and low-energy thread dispatching and synchronization
- Automatic clock gating of PEs waiting on a barrier
- Event-based computation (between accelerators/DMAs and GP cores)

Angelo Garofalo

## Tightly-Coupled Acceleration: Success Stories



Angelo Garofalo

## GAP9 with HWPE (NE16)

		I .	1	1			Benchmark Results								
				Accelerator(s) & Number			Task	Visual Wake Words Visual Wake Words Dataset MobileNetV1 (0.25x) 80% (top 1)		Image Classification CIFAR-10 ResNet-V1 85% (top 1)		n Keyword Spotting Google Speech Commands DSCNN 90% (top 1)		Anomaly Detection ToyADMOS (ToyCar) FC AutoEncoder 0.85 (AUC)	
							Data								
		1					Model Accuracy Units								
		1	Processor(s) & Number		Software										
Submitter	Board Name	SoC Name				Notes		Latency in	Energy in	Latency in E	Energy in	Latency in Energy in	Latency in F	Energy in	
								ms	uJ	ms	uJ	ms	uJ	ms	uJ
Greenwaves			RISC-V Core			GAP9 (370MHZ.	-				-				
Technologies	GAP9 EVK	GAP9	(1+9)	NE16 (1)	GreenWaves GAPFlow	0.8Vcore)		1.13	58.4	0.62	40.4	4 0.48	26.7	0.18	7.29
Greenwaves Technologies	GAP9 EVK	GAP9	RISC-V Core (1+9)	NE16 (1)	GreenWaves GAPFlow	GAP9 (240MHZ, 0.65Vcore)		1.73	40.8	0.95	27.7	0.73	18.6	0.27	5.25
OctoML	NRF5340DK	nRF5340	Arm® Cortex®- M33		microTVM using CMSIS-NN backend	128MHz		232.0		316.1		76.1		6.27	
OctoML	NUCLEO-L4R5ZI	STM32L4R5Z IT6U	Arm® Cortex®- M4		microTVM using CMSIS-NN backend	120MHz, 1.8Vbat		301.2	15531.4	389.5	20236.3	3 99.8	5230.3	8.60	443.2
o	NUCLEO-L4R5ZI	STM32L4R5Z IT6U	Arm® Cortex®- M4		microTVM using native contegen	120MHz, 1.8∨bat		336.5	17131.6	389.2	21342.5	144.0	7950.5	11.7	633.7
OctoML				4					2		8-			10.00	
Plumerai	B_U585I_IOT02A	STM32U585	Arm® Cortex®- M33												
Plumerai Plumerai	B_U585I_IOT02A CY8CPROTO-062- 4343w	STM32U585 PSoC 62 MCU	Arm® Cortex®- M33 Arm® Co M4	Work	ing at best E.I				Мо	bile	Net	tV1	inf	erer	nce
Plumerai Plumerai Plumerai	B_U585I_IOT02A CY8CPROTO-062 4343w DISCO-F746NG	STM32U585 PSoC 62 MCU STM32F746	Arm® Cortex®- M33 Arm® Co M4 Arm® Co M7	Work	king at best E.I				Мо	bile	Net	tV1	info	erer	ice
Plumeral Plumeral Plumeral Plumeral	B_U585I_IOT02A CY8CPROTO-062- 4343w DISCO-F746NG NUCLEO-L4R5ZI	STM32U585 PSoC 62 MCU STM32F746 STM32L4R5Z IT6U	Arm® Cortex®- M33 Arm® Co M4 Arm® Co M7 Arm® Co M4	Work op	king at best E.I erating point				Mo i	bile n 1	Ne .73	tV1 ms	inf at 4	erer 1u]	ice
OctoML Plumerai Plumerai Plumerai Silicon Labs	B_U5851_IOT02A CY8CPROTO-062- 4343w DISCO-F746NG NUCLEO-L4R5ZI xG24-DK2601B	STM32U585 - PSoC 62 MCU STM32F746 STM32L4R5Z IT6U EFR32MG24	Arm® Contexe- M33 Arm® Co M4 Arm® Co M7 Arm® Co M4 Arm® Cortex®- M33	Work op Silicon Labs MVP(1)	ting at best E.I erating point	<b>■</b>		111.6	Mo 1139.2	bile n 1	Net .73	tV1 ms	infe at 4	erer 1u]	1Ce
Plumerai Plumerai Plumerai Plumerai Silicon Labs STMicroelectronics	B_U585I_IOT02A CY8CPROTO-062 4343w DISCO-F746NG NUCLEO-L4R5ZI xG24-DK2601B NUCLEO-H7A3ZI- Q	STM32U585 PSoC 62 MCU STM32F746 STM32L4R5Z IT6U EFR32MG24 STM32H7A3Z IT6Q	Arm® Contexes- M33 Arm® C M4 Arm® Co M7 Arm® Contex®- M33 Arm® Contex®- M33 Arm® Contex®- M7	Work op Silicon Labs MVP(1)	ting at best E.I erating point TensorFlowLite for Microcontrollers, CMSIS-NN, Silicon Labs Gecko SDK X-CUBE-AI v7.3.0	280MHz, 3.3Vbat		111.6	<b>Mo</b> 1139.2 7978.5	bile n 1 120.9 54.3	2Net .73	tV1 ms 7 36.3 3 16.8	inf( at 4 401.5 2721.8	erer 1u3	ACE
Plumerai Plumerai Plumerai Plumerai Silicon Labs STMicroelectronics STMicroelectronics	B_U585I_IOT02A CY8CPROTO-062 4343w DISCO-F746NG NUCLEO-L4R5ZI xG24-DK2601B NUCLEO-H7A3ZI- Q NUCLEO-L4R5ZI	STM32U585 PSoC 62 MCU STM32F746 STM32L4R5Z IT6U EFR32MG24 STM32L4R5Z IT6Q STM32L4R5Z IT6U	Arm® Contexes- M33 Arm® Co M4 Arm® Co M7 Arm® Contex®- M33 Arm® Contex®- M3 Arm® Contex®- M7 Arm® Contex®- M7 Arm® Contex®- M7	Work op Silicon Labs MVP(1)	ting at best E.I erating point TensorFlowLife for Microcontrollers, CMSIS-NN, Silicon Labs Gecko SDK X-CUBE-AI v7.3.0 X-CUBE-AI v7.3.0	280MHz, 3.3Vbat 120MHz, 1.8Vbat		111.6 50.7 230.5	<b>Mo</b> 1139.2 7978.5 10066.6	bile n 1 120.9 54.3 226.9	2Ne .73	<b>tV1</b> ms 7 36.3 3 16.8 3 75.1	info at 4 401.5 2721.8 3371.7	erer 1u3	47.3 266.5 323.0
Plumerai Plumerai Plumerai Plumerai Silicon Labs STMicroelectronics STMicroelectronics	B_U5851_IOT02A CY8CPROTO-062 4343w DISCO-F746NG NUCLEO-L4R5ZI xG24-DK2601B NUCLEO-H7A3ZI- Q NUCLEO-H7A3ZI- Q NUCLEO-L4R5ZI NUCLEO-U575ZI- Q	STM32U585 PSoC 62 MCU STM32F746 STM32L4R52 IT8U EFR32MG24 STM32H7A32 IT6Q STM32L4R52 IT6Q STM32L4R52 IT6Q STM32U5752 IT6Q	Arm® Contexes- M33 Arm® Cr M4 Arm® Cr M7 Arm® Cortex®- M33 Arm® Cortex®- M7 Arm® Cortex®- M4 Arm® Cortex®- M4 Arm® Cortex®- M4 Arm® Cortex®- M3	Work op Silicon Labs MVP(1)	TensorFlowLite for Microcontrollers, CMSIS-NN, Silicon Labs Gecko SDK X-CUBE-AI v7.3.0 X-CUBE-AI v7.3.0 X-CUBE-AI v7.3.0	280MHz, 3.3Vbat 120MHz, 1.8Vbat 160MHz, 1.8Vbat		111.6 50.7 230.5 133.4	<b>Mo</b> 11139.2 7978.5 10066.6 3364.5	bile n 1 120.9 54.3 226.9 139.7	1234.7 8707.3 10681.6 3642.0	tV1 ms ( 7 36.3 3 16.8 3 75.1 0 44.2	401.5 2721.8 3371.7 1138.5	erer 1u3	47.3 266.5 323.0 119.1
Plumerai Plumerai Plumerai Plumerai Silicon Labs STMicroelectronics STMicroelectronics STMicroelectronics STMicroelectronics Syntiant	B_U585I_IOT02A CY8CPROTO-062 4343w DISCO-F746NG NUCLEO-L4R5ZI xG24-DK2601B NUCLEO-U4R5ZI Q NUCLEO-L4R5ZI NUCLEO-U575ZI- Q NDP9120-EVL	STM32U585 PSoC 62 MCU STM32F746 STM32L4R5Z IT6U EFR32MG24 STM32L4R5Z IT6Q STM32L4R5Z IT6Q STM32L4R5Z IT6Q NDP120	Arm® Contexes- M33 Arm® Cr M4 Arm® Cr M7 Arm® Contex®- M33 Arm® Contex®- M7 Arm® Contex®- M7 Arm® Contex®- M4 Arm® Contex®- M4 Arm® Contex®- M3 M0 + HiFi	Work op Silicon Labs MVP(1)	TensorFlowLite for Microcontrollers, CMSIS-NN, Silicon Labs Gecko SDK X-CUBE-AI v7.3.0 X-CUBE-AI v7.3.0 X-CUBE-AI v7.3.0 X-CUBE-AI v7.3.0	280MHz, 3.3Vbat 120MHz, 1.8Vbat 160MHz, 1.8Vbat Syntiant Core 2 (98MHz, 1	         	111.6 50.7 230.5 133.4 4.10	Mo 1139.2 7978.5 10066.6 3364.5 97.2	bile n 1 120.9 54.3 226.9 139.7 5.12	1234.7 8707.3 10681.6 3642.0 139.4	tV1 ms 4 7 36.3 3 16.8 3 75.1 9 44.2 4 1.48	401.5 2721.8 3371.7 1138.5 43.8	5.43 5.43 1.82 7.57 4.84	47.3 266.5 323.0 119.1
Plumerai Plumerai Plumerai Plumerai Silicon Labs STMicroelectronics STMicroelectronics STMicroelectronics Syntiant Syntiant	B_U585I_IOT02A CY8CPROTO-062 4343w DISCO-F746NG NUCLEO-L4R5ZI xG24-DK2601B NUCLEO-U4R5ZI Q NUCLEO-U4R5ZI NUCLEO-U575ZI- Q NDP9120-EVL NDP9120-EVL	STM32U585 PSoC 62 MCU STM32F746 STM32L4R5Z IT6U STM32L4R5Z IT6Q STM32L4R5Z IT6Q STM32L4R5Z IT6Q STM32L4R5Z IT6Q NDP120 NDP120	Arm® Contexes- M33 Arm® Cr M4 Arm® Cr M7 Arm® Contex®- M33 Arm® Contex®- M33 Arm® Contex®- M33 Arm® Contex®- M4 Arm® Contex®- M3 Arm® Contex®- M3 Arm® Contex®- M4 Arm® Contex®- M3 Arm® Contex®- M4 Arm® Contex®- Arm® Contex®- Arm Arm® Contex®- Arm® Contex®- Arm® Contex®- Ar	Work op Silicon Labs MVP(1) Syntiant Core 2 (98MH Syntiant Core 2 (30MH	TensorFlowLite for Microcontrollers, CMSIS-NN, Silicon Labs Gecko SDK X-CUBE-AI v7.3.0 X-CUBE-AI v7.3.0 X-CUBE-AI v7.3.0 tsyntiant TDK	280MHz, 3.3Vbat 120MHz, 1.8Vbat 160MHz, 1.8Vbat Syntiant Core 2 (98MHz, Syntiant Core 2 (30MHz, 1	1	111.6 50.7 230.5 133.4 4.10 12.7	Mo 1139.2 7978.5 10066.6 3364.5 97.2 71.7	bile n 1 120.9 54.3 226.9 139.7 5.12 16.0	1234.7 8707.3 10681.6 3642.0 139.4 101.8	TV1 7 36.3 3 16.8 3 75.1 0 44.2 4 1.48 3 4.37	401.9 2721.8 3371.7 1138.5 33.5	5.43 5.43 1.82 7.57 4.84	1CE 47.3 266.5 323.0 119.1

Best-in-class in **latency** and **energy efficiency** in MLPerf Tiny 1.0 ...by two-to-three orders of magnitude

Angelo Garofalo

## Heterogeneous In-Memory Computing Cluster



[Garofalo et al. IEEE JETCAS, 2022]

- □ 34 AIMC Xbars, 25mm<sup>2</sup> layout in 22nm FDX technology
- Enough storage capacity to host a full MobileNetV2 (~1MB) execution
- Sequential execution model tailored for the needs of an advanced IoT device

Angelo Garofalo

## Heterogeneous In-Memory Computing Cluster



[Ottavi et al. AICAS, 2021] [Khaddam et al. JSSCC, 2022] □ 256x256 PCM-based IMC array

- Weights precision: 4-bit
- Input/Output precision: 8-bit
- MVM operations: 130K in 130ns
   Peak performance **1TOPS**
- Peak Efficiency 11.7 TOPS/W
- Integrated into HWPE
- Data\_itf sized to sustain BW requirement of the IMC core

Angelo Garofalo

## Roof-Line of an IMC Heterogeneous Cluster

![](_page_31_Figure_1.jpeg)

Compute roof determined by MVM op. latency, XBAR size and array utilization
 BW determined by architecture and physical implementation of digital system
 Pipelined execution (overlap computation with streaming) to fully utilize BW

Angelo Garofalo

### End-to-End MobileNetV2 Execution

![](_page_32_Figure_1.jpeg)

Angelo Garofalo

ISSCC 2023 - Forum 6.4: Is an AI accelerator all you need? Overcoming Amdahl's law with tightly-coupled heterogeneous accelerators. 33 of 44

## Outline

- □ The quest for energy-efficient AI at the Edge of the IoT
  - The performance-power proportionality goal
  - Dealing with the end of the Dennard's scaling
- □ Low-power heterogeneous computing: mitigating deep acceleration effects
  - Multi-core heterogeneous systems
  - Tightly-coupled digital acceleration
  - Heterogeneous analog in-memory computing (AIMC) systems
- □ Insights for scaled-up heterogeneous powerful architectures
  - AIMC many-core architecture
  - Asymmetric chiplet-based systems
- Conclusions & Outlook

Angelo Garofalo

## Explorative Scaled-up AIMC Many-Core System

![](_page_34_Figure_1.jpeg)

512 cluster, Low-Latency Hierarchical NoC;
 Not anymore in the IoT Domain (multi-core AIMC parallelism, batching exec.)

Angelo Garofalo

ISSCC 2023 - Forum 6.4: Is an AI accelerator all you need? Overcoming Amdahl's law with tightly-coupled heterogeneous accelerators. 35 of 44

## Exploring System-Level Inefficiencies (ResNet-18)

![](_page_35_Figure_1.jpeg)

- □ HW-aware design space exploration framework to balance resources
   □ More Heterogeneity → Analog Clusters, Purely Digital, Purely SW, Mem..
- Custom execution scheduling to increase HW resource utilization
- On-Chip communication seems not a big issue..
  - □ ..but strong assumptions made on chip area: 480mm<sup>2</sup> → Chiplets!

## Leveraging Chiplet Technology: Occamy example

![](_page_36_Figure_1.jpeg)

- Mitigate slow-down of Moore's Law and end of Dennard's Scaling
- Enable heterogeneity at chiplet level (also different technology nodes)
   Unified chiplet interfaces required to ease integration & packaging

Angelo Garofalo

## Occamy Chiplet

![](_page_37_Figure_1.jpeg)

- □ 216 "snitch" cores/chiplet, **384** GDFLOPS/chiplet @ 1 GHz
- □ On-chip memory: 8-ch 16 GB HBM2e @ 512 GB/s, 2MB + 512kB per cluster SPM
- □ Hierarchical on-chip interconnect based on AXI4 protocol
- □ D2D link: 304 Gb/s @ 250 MHz (duplex)

Angelo Garofalo

ISSCC 2023 - Forum 6.4: Is an AI accelerator all you need? Overcoming Amdahl's law with tightly-coupled heterogeneous accelerators. 38 of 44

## Conclusions & Outlook

- Efficient data movements and mitigation of Amdahl's effects are key challenges to address to scale up AI efficiency on low-power but also high-performance systems
- Current trend towards multi-core heterogeneous AI systems
  - Many heterogeneous accelerators tightly-coupled to SW cores
  - Fast and efficient on-chip communication
  - Event-driven execution modes for on-demand resource scheduling
- Challenges and Outlook
  - System-level integration of multi-core AIMC processors
  - Heterogeneous hardware design space exploration
  - Leveraging chiplet technology to deploy massive heterogeneity in a single system
  - Optimized hardware-aware AI schedulers & compilers to improve data reuse (i.e. reduce data movements), HW utilization, overall system compute efficiency

Angelo Garofalo

## References

- 1. Bianco, Simone, Remi Cadene, Luigi Celona, and Paolo Napoletano. "Benchmark analysis of representative deep neural network architectures." IEEE Access 6 (2018): 64270-64277.
- 2. Rusci, M., Capotondi, A., & Benini, L. (2020). Memory-driven mixed low precision quantization for enabling deep network inference on microcontrollers. *Proceedings of Machine Learning and Systems*, *2*, 326-335.
- 3. Ottavi, G., Garofalo, A., Tagliavini, G., Conti, F., Benini, L., & Rossi, D. (2020, July). A mixed-precision RISC-V processor for extreme-edge DNN inference. In 2020 IEEE Computer Society Annual Symposium on VLSI (ISVLSI) (pp. 512-517). IEEE.
- 4. Garofalo, A., Ottavi, G., Di Mauro, A., Conti, F., Tagliavini, G., Benini, L., & Rossi, D. (2021, September). 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. In ESSCIRC 2021-IEEE 47th European Solid State Circuits Conference (ESSCIRC) (pp. 267-270). IEEE.
- 5. Garofalo, A., Ottavi, G., Conti, F., Karunaratne, G., Boybat, I., Benini, L., & Rossi, D. (2022). A Heterogeneous In-Memory Computing Cluster For Flexible End-to-End Inference of Real-World Deep Neural Networks. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*.
- 6. Bruschi, N., Tagliavini, G., Garofalo, A., Conti, F., Boybat, I., Benini, L., & Rossi, D. (2022). End-to-End DNN Inference on a Massively Parallel Analog In Memory Computing Architecture. *arXiv preprint arXiv:2211.12877*.
- 7. Garofalo, A., Tortorella, Y., Perotti, M., Valente, L., Nadalini, A., Benini, L., ... & Conti, F. (2022). Darkside: A Heterogeneous RISC-V Compute Cluster for Extreme-Edge On-Chip DNN Inference and Training. *IEEE Open Journal of the Solid-State Circuits Society*.
- 8. Gautschi, M., Schiavone, P. D., Traber, A., Loi, I., Pullini, A., Rossi, D., ... & Benini, L. (2017). Near-threshold RISC-V core with DSP extensions for scalable IoT endpoint devices. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 25(10), 2700-2713.
- 9. Ottavi, G., Karunaratne, G., Conti, F., Boybat, I., Benini, L., & Rossi, D. (2021, June). End-to-end 100-TOPS/W Inference With Analog In-Memory Computing: Are We There Yet?. In 2021 IEEE 3rd International Conference on Artificial Intelligence Circuits and Systems (AICAS) (pp. 1-4). IEEE.

Angelo Garofalo

# References - 2

- Glaser, F., Tagliavini, G., Rossi, D., Haugou, G., Huang, Q., & Benini, L. (2020). Energy-efficient hardware-accelerated synchronization for shared-L1-memory multiprocessor clusters. IEEE Transactions on Parallel and Distributed Systems, 32(3), 633-648.
- 11. Zaruba, F., Schuiki, F., Hoefler, T., & Benini, L. (2020). Snitch: A tiny pseudo dual-issue processor for area and energy efficient execution of floating-point intensive workloads. IEEE Transactions on Computers, 70(11), 1845-1860.
- 12. Zaruba, F., Schuiki, F., & Benini, L. (2020). Manticore: A 4096-core RISC-V chiplet architecture for ultraefficient floating-point computing. IEEE Micro, 41(2), 36-42.
- 13. Rossi, D., Conti, F., Eggiman, M., Di Mauro, A., Tagliavini, G., Mach, S., ... & Benini, L. (2021). Vega: A Ten-Core SoC for IoT Endnodes With DNN Acceleration and Cognitive Wake-Up From MRAM-Based State-Retentive Sleep Mode. IEEE Journal of Solid-State Circuits, 57(1), 127-139.
- 14. Khaddam-Aljameh, R., Stanisavljevic, M., Mas, J. F., Karunaratne, G., Brändli, M., Liu, F., ... & Eleftheriou, E. (2022). HERMEScore—a 1.59-TOPS/mm 2 PCM on 14-nm CMOS in-memory compute core using 300-ps/LSB linearized CCO-based ADCs. IEEE Journal of Solid-State Circuits, 57(4), 1027-1038.
- 15. Fick, L., Skrzyniarz, S., Parikh, M., Henry, M. B., & Fick, D. (2022, February). Analog Matrix Processor for Edge AI Real-Time Video Analytics. In 2022 IEEE International Solid-State Circuits Conference (ISSCC) (Vol. 65, pp. 260-262). IEEE.
- 16. Jia, H., Ozatay, M., Tang, Y., Valavi, H., Pathak, R., Lee, J., & Verma, N. (2021, February). 15.1 a programmable neural-network inference accelerator based on scalable in-memory computing. In 2021 IEEE International Solid-State Circuits Conference (ISSCC) (Vol. 64, pp. 236-238). IEEE.
- Ueyoshi, K., Papistas, I. A., Houshmand, P., Sarda, G. M., Jain, V., Shi, M., ... & Verhelst, M. (2022, February). DIANA: An End-to-End Energy-Efficient Digital and ANAlog Hybrid Neural Network SoC. In 2022 IEEE International Solid-State Circuits Conference (ISSCC) (Vol. 65, pp. 1-3). IEEE.
- 18. Jia, H., Valavi, H., Tang, Y., Zhang, J., & Verma, N. (2020). A programmable heterogeneous microprocessor based on bitscalable in-memory computing. IEEE Journal of Solid-State Circuits, 55(9), 2609-2621.

Angelo Garofalo

# References - 3

- 19. Rossi, D., Pullini, A., Loi, I., Gautschi, M., Gürkaynak, F. K., Teman, A., ... & Benini, L. (2017). Energy-efficient near-threshold parallel computing: The PULPv2 cluster. Ieee Micro, 37(5), 20-31.
- Moons, B., Uytterhoeven, R., Dehaene, W., & Verhelst, M. (2017, February). 14.5 envision: A 0.26-to-10tops/w subword-parallel dynamic-voltage-accuracy-frequency-scalable convolutional neural network processor in 28nm fdsoi. In 2017 IEEE International Solid-State Circuits Conference (ISSCC) (pp. 246-247). IEEE.
- 21. Chen, Y. H., Krishna, T., Emer, J. S., & Sze, V. (2016). Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks. IEEE journal of solid-state circuits, 52(1), 127-138.
- 22. Ottavi, G., Garofalo, A., Tagliavini, G., Conti, F., Di Mauro, A., Benini, L., & Rossi, D. (2022). Dustin: A 16-Cores Parallel Ultra-Low-Power Cluster with 2b-to-32b Fully Flexible Bit-Precision and Vector Lockstep Execution Mode. arXiv preprint arXiv:2201.08656.
- 23. Bruschi, N., Tagliavini, G., Garofalo, A., Conti, F., Boybat, I., Benini, L., & Rossi, D. (2022). End-to-End DNN Inference on a Massively Parallel Analog In Memory Computing Architecture. arXiv preprint arXiv:2211.12877. [To appear at DATE 2023].
- 24. Burrello, A., Garofalo, A., Bruschi, N., Tagliavini, G., Rossi, D., & Conti, F. (2021). Dory: Automatic end-to-end deployment of real-world dnns on low-cost iot mcus. IEEE Transactions on Computers, 70(8), 1253-1268.
- 25. Di Mauro, A., Scherer, M., Rossi, D., & Benini, L. (2022). Kraken: A direct event/frame-based multi-sensor fusion soc for ultraefficient visual processing in nano-uavs. arXiv preprint arXiv:2209.01065.
- 26. Jain, V., Giraldo, S., De Roose, J., Boons, B., Mei, L., & Verhelst, M. (2022, June). TinyVers: A 0.8-17 TOPS/W, 1.7 μW-20 mW, Tiny Versatile System-on-chip with State-Retentive eMRAM for Machine Learning Inference at the Extreme Edge. In 2022 IEEE Symposium on VLSI Technology and Circuits (VLSI Technology and Circuits) (pp. 20-21). IEEE.

Angelo Garofalo

#### Is an AI Accelerator All You Need? Overcoming Amdahl's Law With Tightly-Coupled Specialized Accelerators

**BACK-UP SLIDES** 

ISSCC 2023, FORUM 6 "The Future of Heterogeneous Multi-Core Architectures for Al and Other Specialized Processing"

## From High-Performance to Low-Power

- Multi-Core Architectures
  - Speed-up general-purpose tasks
  - Domain-specific insns rather than super-scalar cores for higher en. efficiency

Dedicated acceleration solutions to boost AI workloads

- Processor specialization (low-bitwidth integer arithmetic)
- Dedicated digital/analog data-paths

□ Heterogeneous Integration to mitigate deep acceleration effects

- CPUs + domain-specific GP cores + NN Engines
- High-Bandwidth Low-Latency On-Chip Communication

Angelo Garofalo

### The Tunnel: High-Performance vs. High-Efficiency

![](_page_44_Figure_1.jpeg)

Angelo Garofalo

#### Extended Dot-Product Unit

![](_page_45_Figure_1.jpeg)

## Mixed-Precision Controller

![](_page_46_Figure_1.jpeg)

Mixed-Precision operations
require a controller:
selection of the correct subword of the lowest precision
operand (Vector B) to be
used in current SIMD op.

The controller is programmed by control status registers (CSRs).

Angelo Garofalo

## **Dynamic Bit-Scalable Execution**

Standard Instructions							Virtua	al <b>Mistua</b>	tionstructions	
pv.dotsp.h pv.dotsp.b pv.dotsp.n pv.dotsp.max2 pv.dotsp.m8x2 pv.dotsp.m8x4 pv.dotsp.m16x8 pv.dotsp.m16x8 pv.dotsp.m16x4 pv.dotsp.m16x2 pv.dotsp.sc.h pv.dotsp.sc.c pv.dotsp.sc.n pv.dotsp.sc.n pv.dotsp.sc.n pv.dotsp.sc.m8x2 pv.dotsp.sc.m8x4 pv.dotsp.sc.m16x8 pv.dotsp.sc.m16x8 pv.dotsp.sc.m16x8 pv.dotsp.sc.m16x8 pv.dotsp.sc.m16x2 pv.dotsp.sc.m16x2	pv.dotup.h pv.dotup.b pv.dotup.n pv.dotup.r pv.dotup.m4x2 pv.dotup.m8x2 pv.dotup.m8x4 pv.dotup.m16x8 pv.dotup.m16x8 pv.dotup.m16x2 pv.dotup.sc.h pv.dotup.sc.h pv.dotup.sc.n pv.dotup.sc.n pv.dotup.sc.n pv.dotup.sc.m8x2 pv.dotup.sc.m8x4 pv.dotup.sc.m16x8 pv.dotup.sc.m16x8 pv.dotup.sc.m16x8 pv.dotup.sc.m16x8 pv.dotup.sc.m16x2 pv.dotup.sc.ih pv.dotup.sc.ih	pv.dotusp.h pv.dotusp.n pv.dotusp.n pv.dotusp.m4x2 pv.dotusp.m8x2 pv.dotusp.m8x2 pv.dotusp.m16x8 pv.dotusp.m16x8 pv.dotusp.m16x4 pv.dotusp.m16x2 pv.dotusp.sc.h pv.dotusp.sc.h pv.dotusp.sc.n pv.dotusp.sc.n pv.dotusp.sc.m8x2 pv.dotusp.sc.m16x8 pv.dotusp.sc.m16x8 pv.dotusp.sc.m16x4 pv.dotusp.sc.m16x4 pv.dotusp.sc.m16x4 pv.dotusp.sc.m16x2 pv.dotusp.sc.m16x2 pv.dotusp.sc.m16x2	pv.sdotsp.h pv.sdotsp.b pv.sdotsp.n pv.sdotsp.c pv.sdotsp.m8x2 pv.sdotsp.m8x2 pv.sdotsp.m8x4 pv.sdotsp.m16x8 pv.sdotsp.m16x4 pv.sdotsp.m16x4 pv.sdotsp.sc.h pv.sdotsp.sc.c pv.sdotsp.sc.c pv.sdotsp.sc.n pv.sdotsp.sc.m8x2 pv.sdotsp.sc.m8x4 pv.sdotsp.sc.m16x8 pv.sdotsp.sc.m16x8 pv.sdotsp.sc.m8x4 pv.sdotsp.sc.m16x2 pv.sdotsp.sc.m16x2 pv.sdotsp.sc.m16x2 pv.sdotsp.sc.m16x2 pv.sdotsp.sc.m16x2 pv.sdotsp.sc.m16x2	pv:sdotup.h pv:sdotup.n pv:sdotup.n pv:sdotup.m pv:sdotup.m&x2 pv:sdotup.m&x2 pv:sdotup.m8x4 pv:sdotup.m16x8 pv:sdotup.m16x4 pv:sdotup.m16x2 pv:sdotup.sc.h pv:sdotup.sc.c pv:sdotup.sc.c pv:sdotup.sc.n pv:sdotup.sc.m8x2 pv:sdotup.sc.m8x4 pv:sdotup.sc.m16x8 pv:sdotup.sc.m16x8 pv:sdotup.sc.m16x4 pv:sdotup.sc.m16x4 pv:sdotup.sc.h pv:sdotup.sc.h pv:sdotup.sc.h	pv.sdotusp.h pv.sdotusp.b pv.sdotusp.c pv.sdotusp.c pv.sdotusp.m4x2 pv.sdotusp.m8x2 pv.sdotusp.m16x4 pv.sdotusp.m16x4 pv.sdotusp.sc.h pv.sdotusp.sc.b pv.sdotusp.sc.c pv.sdotusp.sc.c pv.sdotusp.sc.m8x2 pv.sdotusp.sc.m8x4 pv.sdotusp.sc.m8x4 pv.sdotusp.sc.m8x4 pv.sdotusp.sc.m16x8 pv.sdotusp.sc.m16x8 pv.sdotusp.sc.m16x8 pv.sdotusp.sc.m16x4 pv.sdotusp.sc.m16x4	] ] ] [e 10	• No e Reus c Bit- e Exec	pv.dotsp.v pv.dotsp.sc .v pv.dotsp.sc .v pv.dotup.sc .v pv.dotup.sc .v pv.dotup.sc .v pv.dotusp.sc pv.dotusp.sc	<pre>pv.sdotsp.v pv.sdotsp.sc.v pv.sdotsp.sci.v pv.sdotup.sci.v pv.sdotup.sci.v pv.sdotup.sci.v pv.sdotusp.sci.v pv.sdotusp.sci.v pv.sdotusp.sci.v int main() { </pre>	ed at ID insn ats <b>bde</b>
pv.dotsp.sci.c pv.dotsp.sci.m4x2 pv.dotsp.sci.m8x2 pv.dotsp.sci.m8x4 pv.dotsp.sci.m16x8 pv.dotsp.sci.m16x4 pv.dotsp.sci.m16x4	pv.dotup.sci.c pv.dotup.sci.m4x2 pv.dotup.sci.m8x2 pv.dotup.sci.m8x4 pv.dotup.sci.m16x8 pv.dotup.sci.m16x4 pv.dotup.sci.m16x2	pv.dotusp.sci.c pv.dotusp.sci.n pv.dotusp.sci.m4x2 pv.dotusp.sci.m8x2 pv.dotusp.sci.m8x4 pv.dotusp.sci.m16x8 pv.dotusp.sci.m16x4 pv.dotusp.sci.m16x2	pv.sdotsp.sci.c pv.sdotsp.sci.n pv.sdotsp.sci.m4x2 pv.sdotsp.sci.m8x2 pv.sdotsp.sci.m8x4 pv.sdotsp.sci.m16x8 pv.sdotsp.sci.m16x2 instru	pv.sdotup.sci.c pv.sdotup.sci.m4x2 pv.sdotup.sci.m8x2 pv.sdotup.sci.m8x4 pv.sdotup.sci.m16x8 pv.sdotup.sci.m16x4 pv.sdotup.sci.m16x2	pv:sdotusp.sci.c pv:sdotusp.sci.m pv:sdotusp.sci.m4x2 pv:sdotusp.sci.m8x4 pv:sdotusp.sci.m8x4 pv:sdotusp.sci.m16x8 pv:sdotusp.sci.m16x2	p. ne inst	v x20,>	, <11,x10	SIMD_FMT( convolution(  SIMD_FMT( convolution(	M8x4); [A, W, Res); M8x2); [A, W, Res);
pv.packl pv.sdots	ni.b x15, x sp.b x20, x	7, x8 15, x10	per fo and t	ormat ype		per	type		}	

Angelo Garofalo

## Mac-Load Instruction

![](_page_48_Figure_1.jpeg)

□ *MatMul's operands* memory access patern extremely regular

- And known before execution
- □ Fuse MAC with LOAD operations in one single-cycle insns
  - "Fused" data-path, contemporary cosume current operand & prefetch next-one
  - Reduce overhead of Loads into MatMul innermost loops

Angelo Garofalo

#### Mac-Load: MatMul Innermost Kernel

8-bit innermost loop of the 4x2 Matmul kernel of the PULP-NN library

8-bit innermost loop of the 4x2 MatMul kernel of the extended PULP-NN library with nnsdotp

lp.setup	11, 12,end	HW LOOP		pv.nnsdotusp.h	zero, aw1,16	
p.lw p.lw	W1, 4(aW1!) W2, 4(aW2!)		INIT	pv.nnsdotusp.h	zero, aw3,20	
p.lw	w3, 4(aw3!)	LD/ST WITH POST	NN-RF	pv.nnsdotusp.h	zero, aw4,22	
p.lw	w4, 4(aw4!)	INCREMENT		pv.nnsdotusp.h	zero, ax1,8	
p.lw	x1, 4(ax1!)			lp.setup	11, 12, end	
p.lw	x2, 4(ax2!)			<pre>pv.nnsdotup.h</pre>	zero,ax2,9	
pv.sdotusp.b	s1, x1, w1			pv.nnsdotusp.b	s1, aw2, 0	
pv.sdotusp.b	s2, x1, w2	8-B SIMD MAC		<pre>pv.nnsdotusp.b</pre>	s2, aw4, 2	
pv.sdotusp.b	s3, x1, w3			<pre>pv.nnsdotusp.b</pre>	s3, aw3, 4	
pv.sdotusp.b	s4, x1, w4			pv.nnsdotusp.b	s4, ax1, 14	
pv.sdotusp.b	s5, x2, w1			<pre>pv.nnsdotusp.b</pre>	s5, aw2, 17	
pv.sdotusp.b	s6, x2, w2	8 SIMD MACs		<pre>pv.nnsdotusp.b</pre>	s6, aw4, 19	8 SIMD MACs
pv.sdotusp.b	s7, x2, w3	with 6		<pre>pv.nnsdotusp.b</pre>	s7, aw3, 21	with 1
<pre>end: pv.sdotusp.b</pre>	s8, x2, w4	explicit LOADs	end:	pv.nnsdotusp.b	s8, aw1, 23	explicit LOAD
Wi Xi : weight/a	ctivation el	Lements Si	: accumul	ators		_

AW<sub>i</sub> AX<sub>i</sub> : addresses for the MEM access

accumulators

L<sub>i</sub> : loop setup

Angelo Garofalo

#### Mac-Load: MatMul Data-Reuse

![](_page_50_Figure_1.jpeg)

#### Up to 94% of SIMD Dotp Unit Utilization on MatMul kernels

Angelo Garofalo

## Working with Unreliability: Quentin (ULP)

![](_page_51_Figure_1.jpeg)

Angelo Garofalo

## Working with Unreliability: Quentin (ULP)

![](_page_52_Figure_1.jpeg)

**Perf Tuning:** 22mW@140Gop/s 635µW@20Gop/s

Angelo Garofalo

ISSCC 2023 - Forum 6.4: Is an AI accelerator all you need? Overcoming Amdahl's law with tightly-coupled heterogeneous accelerators.

#### End-to-End MobileNetV2 on Darkside

![](_page_53_Figure_1.jpeg)

- □ Heterogeneous acceleration mitigates Amdahl's effects on *Bottleneck*s
- Still room to accelerate PW convolutions

Angelo Garofalo

## MobileNetV2 mapping on AIMC

![](_page_54_Figure_1.jpeg)

Only PointWise layers mapped on AIMC

Tile&Pack algorithm to deploy MBNTv2 on the smallest number of AIMC cores
 High storage utilization, poor bit-cell compute utilization on small layers

Angelo Garofalo

## Depth-Wise Execution on AIMC

![](_page_55_Figure_1.jpeg)

- Representative of modern CNN
- Bottlenecks are building block of MobileNet
- DW are used to reduce number of parameters in CNN for mobile applications
- DWs do not map efficiently on IMA

Angelo Garofalo

ISSCC 2023 - Forum 6.4: Is an AI accelerator all you need? Overcoming Amdahl's law with tightly-coupled heterogeneous accelerators.

#### **Execution Mappings**

- □ Full SW solution (8-cores)
- All layers in IMA: 16 channels/job 54% IMA area overhead
- All layers in IMA: 8 channels/job 25% IMA area overhead
- Conv 1x1 on IMA while DW on 8 cores: 0% IMA area overhead
- Conv 1x1 on IMA, DW on Dig Acc., Final Add on SW Cores

## Mitigation of Amdahl's Law on Het. AIMC System

![](_page_56_Figure_1.jpeg)

Angelo Garofalo

## AIMC Performance, Energy and Area Efficiency

#### Execution of Bottleneck layer

![](_page_57_Figure_2.jpeg)

Angelo Garofalo

## AIMC-Based System vs Purely Digital Systems

#### End-to-End execution of MobileNetV2

![](_page_58_Figure_2.jpeg)

- □ Single-core MCU is inefficient to complement AIMC computing capabilities
- □ AIMC heterogeneous cluster for highest performance

Angelo Garofalo

#### AIMC Many-Core Computational Model

![](_page_59_Figure_1.jpeg)

Angelo Garofalo

### AIMC Many-Core System: Balancing The Pipeline

![](_page_60_Figure_1.jpeg)

#### Angelo Garofalo

## Occamy NoC: Efficient and Flexible Data Movement

![](_page_61_Figure_1.jpeg)

DORY: [Burrello et al. IEEE TC, 2021]

- Problem: HBM accesses are critical in terms of
  - Access energy
  - Congestion
  - □ High latency
- Solution: reuse data on lower levels of memory hierarchy
  - Between clusters
  - Across groups
- Smartly distribute workload
  - Clusters: DORY framework for tiling strategies
  - □ Chiplets: layer pipelining, ..

Angelo Garofalo