ISSCC 2024 Forum F1: Efficient Chiplets and Die-to-Die Communications

1.5: Efficient Domain-Specific Compute with Chiplets



Prof. Dejan Marković UCLA ECE Department dejan@ucla.edu

Evolving Standards: Flexibility & Efficiency

Objectives: lower development cost and shorter time-to-market



SoC/ASIC revision / iteration is \$\$\$ (~\$100M in 16nm CMOS)
 Long design cycles (>1 yr)^[3] with increasing design complexity

SoCs Today = CPU/GPU + Accelerators



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

Optimize for Efficiency and Flexibility

Two ways to think about it

Add flexibility to accelerators

Narrow coverage of DSPs

The how

- □ Interconnect
- Switch-boxes
- Sw toolchain



Efficient Multi-Chip Module (MCM) Scaling



Challenges with MCM Design

Challenges:

- High bandwidth density
- Low link latency
- Low energy transfer
- Low I/O area

Chiplet size:

- □ Sweet spot: $\sim 100 \text{ mm}^2$
- UDSP prototype (\$ limited):
 ~6mm²



Research Aims

Domain-specific hardware acceleration

- ASIC-like energy efficiency and throughput
- Just-enough flexibility for a domain
- Key: flexible cores, efficient interconnect

Tile-able chiplets on Silicon Interconnect Fabric (Si-IF)

- Develop scalable interconnects
- Near-range I/O and PHY for cutting-edge bandwidth/latency/energy
- Low-area, portable timing correction circuits for Si-IF I/Os

A 16nm 2x2 Chiplet with 10- μ m Pitch-I/O

Universal Digital Signal Processor (UDSP) Array



[9] U. Rathore, S. Nagi, S. Iyer, D. Markovic, ISSCC 2022.

UDSP Multi-Chip, Multi-Program Tenancy



UDSP Overview



Evolution of UDSP Core



Efficiency and Flexibility in Comm. DSP



Algorithm Ontology: Example DSP Kernels

Example DSP kernels derived from common DSP algorithms

- Up/Down Conv.
- MIMO
- □ IFFF/FFT
- Neural Network
- Zero Forcing
- □ MMSE
- Vector-dot product
- □ MAC, FIR, Euclidian



Iterative Process of Core Design



Balancing core granularity and core utilization to maximize energy and area efficiency

Interconnects: An Exercise in Co-Design



Layer-1 Interconnect (Distance = 1)



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

Layer-2 Interconnect (Distance = $\sqrt{2}$)



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

Layer-3 Interconnect (Distance = 2)



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

Layer-4 Interconnect



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

Layer-4 Interconnect



N-Layer Switch Box: Hyper-Matrix Model

□ Hyper-vector cross-correlation (HVCC) in each dimension (layer)

HVCC for a layer measures inter-dependencies of paths



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

DSE: Search Space Traversal

MCBF & MCBF/HWC plotted against layer density for 3-layer SB



DSE: Maximizing Silicon Area Efficiency

□ Sw/Hw balance is at the peak of MCFB/HWC



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

Si-IF Assembly Overview



Si-IF Characteristics



Si-IF Process and Assembly

□ UDSP dielet powered on to verify Clk tree and shift-registers

- Low-freq Clk applied using a probe station
- □ Dice defect-free Si-IF sites for assembly
 - Template-based wafer scan for repeated patterns
- Dielets assembled on Si-IF using direct Cu-Cu TCB
 - In-situ formic acid (FA) vapor treatment
- □ Ionizers on the bonding tool to ensure an ESD-safe assembly
 - Default 20 I/min flow interferes with the FA vapor flow of 4.5 I/min
 - Leads to inadequate cleaning of Cu pads, inferior bonding quality
 - □ Shear strength <20N, below MIL-SPEC 883G of **50N** for 2mm x 2mm dielets
 - Solution: the ionizer flow rate is optimized w.r.t. FA vapor time

Bonding Parameters of UDSP Dielets



Bonding parameters

- Temperature
- Bonding force
- Bond-head position

Key insights

- FA flow of 4 l/min during placement: bond strength 80N
- No need to shut down ionizer (no ESD events on all 7,168 links tested)

ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

SNR-10 I/O Channel



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

Challenges in Si-IF Link Design

- □ Available area per I/O
 - SerDes^[14] \rightarrow 10,175 μ m²
 - Interposer^[15] \rightarrow 500 μ m²
 - Si-IF/SNR-10 \rightarrow 137 μ m²
- Required throughput
 1.1 Gbps/pin
- Logical redundancy
 Mfg and integration defects

[14] M. Lin, T. Huang, JSSC, 2020. [15] J. M. Wilson, ISSCC, 2018.



Relative per I/O Area

Streaming Near Range (SNR) Channel

- Uni-directional I/Os
- Clock-forwarded
- 64-bits / channel
- Amortized Clk correction and transfer
- Redundancy
- Minimal ESD
- Minimal handshake
- Sync/Async modes
- 3 Clk cycles of latency
- FIFO for CDT



Tx/Rx Cells



Tx pads use minimal additional ESD (64x drive buffer diodes)
 Rx pads use relatively larger ESD protection diodes

Redundancy and Repair: Check



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

Redundancy and Repair: Detect



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

Redundancy and Repair: Repair



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

SNR-10 vs. State-of-the-Art

Reference	[This work]	CICC'21	JSSC'20	ISSCC'18
Techology	16nm	16nm	7nm	16nm
Package Substrate (# layers)	Si-IF (2)	EMIB (4)	CoWoS (15)	MCM/PCB
Bump Pitch (µm)	10	55	40	150
Reach (µm)	350	3,000	500	80,000
Data Rate (Gbps/pin)	1.1	2	8	25
Voltage (V)	0.8	0.9	0.3	0.95
Energy Effc. (pJ/bit)	0.38	0.83	0.56	1.17
I/O Area Density (µm ² /bit)	137	203	500	10,175
Max Shoreline BW (Gbps/mm)	297	256	1600	292
Layer BW (Gbps/mm/layer)	149	64	107	25

SNR-10 Transfer Efficiency vs. Voltage



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

UDSP Control Logic



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

UDSP Compile Flow



UDSP Dynamic Reconfiguration



UDSP Compile Time Breakdown



UDSP 2x2 Test Setup



Si-IF Assembly on PCB



- □ No internal ESD events due to Si-IF assembly
 - Over 6,000 pads tested across two samples
- □ Wire-bonded Si-IF sample on a daughter board with ESD
 - No ESD events observed across 250 pins tested

SNR-10 Performance Summary

- No internal ESD events due to Si-IF assembly
- □ 7 64-b channels per side
 - 448 bits per side
- □ Fmax: 1.1 GHz

□ Inter-die BW: 493 Gbps

□ Max BW density: 297 Gbps/mm



UDSP Fmax, Power & Efficiency vs. V_{DD}



Bridging the Efficiency Gap



Optimal E x A, Throughput Comparison



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

Prior Design Comparison



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

Average Compile Time Comparison



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

Runtime Reconfigurable Array (RTRA)

- □ What problem does it solve?
 - Data-driven dynamic relocation of compute resources
 - Fast response to environmental dynamics or uncertainty
- □ What does it do that others can't?
 - Fast (sub-µs) reconfiguration of Hw resources
 - Multi-program tenancy with priority and multi-size compile
- □ How does RTRA do it?
 - Online scheduler for multi-program, multi-size, priority
 - Interconnect tailored to a computation domain
- Benefit vs. FPGA DSP?
 - 10x higher compute efficiency, 10x utilization, 3-4x throughput

Hw-in-the-Loop RTRA Feasibility Study

 Characterize utilization, time to repurpose, compute time (config + exe) under multi-program, multi-size, priority mappings

Compare with statically configured UDSP and FPGA



Blind Spectrum Sensing Application



Data-driven attentive processing: detection, extraction, classification

Algorithms from: [18] R. Harjani, et al., IEEE Comm. Magazine, Oct. 2015.

RTRA Dynamically Schedules Hw Array



UDSP array size: 18x18



RTRA array size: 18x18

RTRA Maximizes Area Utilization



RTRA Architecture Building Blocks



RT-Compiler/Sch:

Priority-based dynamic multiprogram scheduling, multi-size compile

- IO network: equal latency to all Pes
- Data RAM: kernels with pre-mapped data
- Multi-Bank Memory
 + x-bar: sustains
 high data throughput

[17]

ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

Scalable RTRA Compute Architecture

- □ Key architecture features
 - 16-bit floating-point
 - Modular for scalability
 - Multi-level scheduler
- ~125mm² baseline
 - CPU-sized, HBM2 interface
 - Tile-able chiplet assembly
- Proof-of-concept
 - Integrated Sw/Hw DSE
 - FPGA emulation
 - Adapt SoC to I/O & area



Peripheral I/O: **9.9x more native BW** than HBM2 (for chiplet assembly) doable in 6-layer 10-µm pitch



Software: Compiler Front-End



Software: Compiler Back-End



Future: Hardware Accelerator as a Service

Dynamically compose and relocate accelerators on-the-fly and handle multiple accelerator requests simultaneously



ISSCC 2024 - Forum 1.5: Efficient Domain-Specific Compute with Chiplets

Advancements in Chiplet Assembly

HBM 2 interface

- 1024 bits
 - □ 8ch x 128b (16ch x 64b)
- 89 b/mm

RTRA interface

- 10,112 bits
 - □ 158 PEs per side
 - □ 64 bits per PE
- 879 b/mm
- A 9.9x gap
 - 149b/mm/layer
 - 6-layer 10-µm Si-IF I/O



Peripheral I/O: **9.9x more native BW** than HBM2 (for chiplet assembly) doable in 6-layer 10- μ m pitch



Acknowledgments

- Dr. Uneeb Rathore
- Dr. Sumeet Nagi
- □ Chenkai (Tim) Ling
- □ Prof. Subramanian Iyer
- □ Krutikesh Sahoo
- Dr. Sina Basir-Kazeruni
- Dr. Siva Chandra Jangam
- DARPA CHIPS, DRBE programs
- Accton, Inc.

References (1/2)

- [1] Xilinx Zynq RFSoC DFE
- [2] 3GPP, Online: <u>https://www.3GPP.org</u>
- [3] J. Cong et al., IEEE Design and Test of Computers, 2011.
- [4] Y.S. Shao, B. Reagen, G.Y. Wei, D. Brooks, IEEE Micro 2015, 35(3), pp. 58-70.
- [5] C.C. Wang, UCLA PhD Thesis, 2013.
- [6] F-L. Yuan, UCLA PhD Thesis, 2014.
- [7] L. Leibo, ACM Comput. Survey 2020.
- [8] Seeds, Bose-Einstein Yield Models, Online: <u>https://www.eesemi.com/test-yield-models.htm</u>
- [9] U. Rathore, S. Nagi, S. Iyer, D. Markovic, ISSCC 2022, pp. 52-54.
- [10] C.C. Wang, F.-L. Yuan, T.-H. Yu, and D. Marković, ISSCC 2014, pp. 460-461.
- [11] F.-L. Yuan, T.-H. Yu, D. Markovic, VLSI 2014, pp. 65-66.
- [12] F.L. Yuan, et al., VLSI 2015, pp. 150-151.
- [13] S. Jangam, UCLA PhD Thesis, 2020.
- [14] M. Lin, T. Huang, JSSC, 2020.
- [15] J. M. Wilson, ISSCC, 2018.
- [16] F.-L. Yuan, C.C. Wang, T.-H. Yu, D. Markovic, IEEE JSSC, Jan. 2015, 50(1), pp. 137-149.
- [17] U. Rathore, UCLA PhD Thesis, 2022.
- [18] R. Harjani, et al., IEEE Comm. Magazine, Oct. 2015, 53(10), pp. 173-181.
- [19] S. Nagi, UCLA PhD Thesis, 2022.

References (2/2)

- [20] AMD on Why Chiplets And Why Now, Online: <u>https://www.nextplatform.com/2021/06/09/amd-on-why-chiplets-and-why-now</u>
- [21] Tesla V100 Price, Online: <u>https://www.microway.com/hpc-tech-tips/nvidia-tesla-v100-price-analysis</u>