Modular Chiplet Opportunities for HPC *Evaluating Chiplets for Post Exascale Supercomputing*

John Shalf

Department Head for Computer Science Lawrence Berkeley National Laboratory

Seattle, Washington

ModSim 201





1

Explosion of Computing Demand: Driving Need for Hyper-exponential Improvement in Performance, Energy Efficiency and Integration



Demand for Computing

Supply for Computing





NVIDIAnomics



 NVIDIA H100 SXM consumes 700W each

- Next Generation B100 is projected to consume 1400W each! (100% increase)
- Street price \$20k-\$30k
 - That's an 800% profit
- And still supply cannot keep up with demand

BERKELEY LAB



Google Sustainability Page in Late 2023

Google Sustainability Empowering individuals Working together Operating sustainably Reports

Q

Overview Net-zero carbon Water stewardship Circular economy Nature & biodiversity Stories

Net-zero carbon



What a difference a year makes!



Google's AI ambitions threaten carbon neutrality pledge

Data center expansion leads to 48% Increase in emissions since 2019

Omer Kabir 09:53, 04.07.24

TAGS: <u>Climate Change</u> <u>Data Center</u> <u>Google</u> <u>Emissions</u>





6



Algorithm-Driven Codesign of Specialized Architectures for Energy-Efficient HPC



NASEM study on post-Exascale computing "We must expand (and create where necessary) integrated teams that identify the key algorithmic and data access motifs in its applications and begin collaborative ab-initio hardware development of supporting accelerators,... a first principles approach that considers alternative mathematical models to account for the limitations of weak scaling."



Neil Thompson: Economics of Post-Moore Electronics

http://neil-t.com, MIT CSAIL, MIT Sloan School



ա....վ

BERKELEY LAB



Why? Domain specific Architectures driven by hyperscalers

in response to slowing of Moore's Law (switch to systems focus for future scaling)

BERKELEY LAB



What is a Chiplet?





Learning from Attack of the Killer Micros



Attack of the killer micros John Markoff, May 6, 1991





Technology Insertion into Mainstream Platforms

AMD, Intel, Arm offer integration path for 3rd party accelerator "chiplets"

Modular AMD Chips to Embrace Custom 3rd To 'Meteor Lake' and Beyond: How **Party Chiplets**

By Francisco Pires last updated June 20, 2022 News

Supercharging learnings - and earnings - from the console space.

🚹 🛯 🚳 😰 🕞 💟 🗭 Comments (2)

When you purchase through links on our site, we may earn an affiliate commission. Here's how it works.

GPU CPU **3rd Party** 6 AIE



Intel Plans a New Era of 'Chiplet'-**Based CPUs**

At the Hot Chips 2022 conference, Intel teased its upcoming 'Meteor Lake' and 'Arrow Lake' processor families, which will use multiple tiny tiles fused together in an attempt to break free of the limits of monolithic chip design. Here's why little tiles are a big deal.







October 19, 2023

It is safe to say that ARM isn't a scrappy startup that was once the pride of the UK. The US-based IPO made the chip designer a big-game chip player, and the new capital is kickstarting some major initiatives to find more customers for its products. A new effort called Total Design aims at making it easier for companies looking to design chips in-house, an idea gaining ground with the Al boom and chip shortages.



Architecture Specialization for Science

(hardware is design around the algorithms) can't design effective hardware without math



Example of Mixed-Radix MultiDimensional FFT Accelerator

Initial steps towards a basis set of hardware accelerator primitives for science



FFT96 Accelerator Die: 16nm TSMC

Area eff.: 4.18 TF/mm² Energy eff.: 4.8 TF/W

NVIDIA H100 in 4nm TSMC (10x

<u>denser than 16nm</u>) Area eff: 0.08 TF/mm² Energy eff: 0.0957 TF/W



Accomplishment: Demonstrated general FFT Accelerator tile generator in 16nm TSMC that outperforms NVIDIA H100 by 50x in raw performance/area and 50x in energy efficiency/flop

-Required only 3 codelet primitives

Goal: Generalize this approach to generators for accelerator hardware primitives that cover broad spectrum of algorithms (*e.g., FFT, Dense/Sparse Linear Algebra, particles and PDEs*)

Mario Vega, Xiaokun Yang, John Shalf, Doru-Thom Popovici: Towards a Flexible Hardware Implementation for Mixed-Radix Fourier Transforms. HPEC 2023: 1-7





The Importance for High Performance DFT

BERKELEY LAB



Circuit/Dataflow Reformulation of DFT



Building Flexible Accelerators

Parameterized Hardware Generation with CHISEL

HH_Core hh_cnt:0~n-1	ſ			
Buffer IIII_datapath br×n St: br				
Timing_ctr				

Table 1.3 Resource Cost Analysis

A Single HH-core

TSQR Designs		FP Operators			Register file (Words)		
Multi-core	Streaming Width	Iterative	Pipeline	Parallel (16 cores)	Iterative	Pipeline	Parallel (16 cores)
TSQR Designs	16	68	432	1,088	192	432	3,072
	32	132	1,632	2,112	768	3672	12,288
	64	260	6,336	4,160	3,072	28,336	49,152
	128	516	24,960	8,256	12,288	222,560	196,608
	256	1,028	99,072	16,448	49,152	1,764,032	786,432



Generate any level of concurrency or pipelining



Modeling the Baseline Computation Limitations of the Offload Model for Accelerators

BERKELEY LAB



Point

iFFT3D

GEMM

Peer Accelerator Model (merged Kernels)

Convert Problem to Compute Bound



ĦF

iFFT1D

Point wise

FFT1D



HPC Microbiome Analysis using GraphBLAS

 Use HPC algorithms and systems for orders of magnitude speedup and to solve previously intractable problems

>>10 ⁰ sequencing reads 36 bp - 1 kb		
	3.65	

Metagenome Assembly	Protein Clustering	Comparative Analysis / Community Detection
Assemble millions of metagenomes based on incomplete data	Cluster billions of proteins	Use fast alignment and annotation for time-sensitive analyses.
Graph algorithms, De-Bruijin graphs, Hash Tables, alignment (Smith-Waterman)	Machine learning (clustering), sparse linear algebra / graphs	Alignment, Machine learning (dimensionality reduction), linear algebra

ASA: Accelerator for GraphBLAS Sparse Accumulation Primitives



Majority of time goes to hypersparse accumulation

- Typically implemented as a hash table
- Replaced has lookup with an associative match
- Reaches maximum efficiency with 16k associative

BERKELEY LAB

Performance Breakdown of HipMCL and SpGEMM



SpGEMM and HipMCL Speedup

HipMCL Timing Breakdown SpGEMM Timing Breakdown 100% Accumulation 100% Waiting 80% bac Buffer 80% Partial Hardware probing 1 Sum 60% and accumulation Probing 2.25x 60% Cache 2.35x Ai * Bj Key 40% 40% 4.55x 5.05x Memory Hierarch Column Merging 20% Address Generator 20% All nonzero entries 0% Keys Tail Pointer Boundary 0% 2 Overflow Values Register Register ASA Baseline ASA-tiling Hardware **Baseline** ASA ASA-tiling Gathering ③ Software check vector size and address Abcast Bcast localspgemm Software sor 4 and merge ■ Col Read ■ Symbolic ■ Sparse Accumulation ■ Sort and Merge Inflation Software Sorting and Conditional Merging multiway-merge Prune Component

One Tile of ASA Architecture Diagram

- ASA accelerator results in 5.05x speedup (average over test inputs) for SpGEMM kernel and 2.35x speedup for HipMCL application overall
- Reduced sparse accumulation cost and *eliminated symbolic computation phase*

Accelerating GraphBLAS Primitives

GraphBLAS¹

Graph algorithms are difficult to implement efficiently in hardware since it requires both graph expertise and hardware expertise

 GraphBLAS separates these concerns by expressing graph algorithms as linear algebraic operations



ASA: Accelerating Sparse Accumulation²

Hardware accelerator that extends the ISA to add dedicated instructions for partial sums and accumulates in sparse matrix multiplication to avoid data dependent branches.

Can this accelerator targeting sparse general matrix-matrix multiplication (SpGEMM) be generalized for a wider range of graph algorithms?

GISA: GraphBLAS Instruction Set Architecture

GraphBLAS to GISA

GraphBLAS Operation	GISA Equivalent
Vector_build Matrix_build	GISA.malloc addr
Matrix_extractTuples	GISA.gather
reduce	GISA.gather_sum
$C_{[:,j]} = A \oplus . \otimes B_{[:,j]}$	for $B_{[k,j]}$ in $B_{[:,j]}$ for $A_{[i,k]}$ in $A_{[:,k]}$ GISA.insert_add i $A_{[i,k]} \otimes B_{[k,j]}$ $C_{[:,j]} = GISA.gather$

The flexibility of the \oplus and \otimes operators enables the expression of a wide range of graph algorithms. GISA captures this by utilizing an ALU to perform the "plus" and "times" operations.

Microarchitecture



Figure 1: Datapath of GISA hardware accelerator. This has been written in SystemVerilog and is parameterizable by size and associativity

BERKELEY LAB

Data Rearrangement Engine (DRE) Near Memory

Xueyang Liu, Patricia Gonzalez-Guerrero, Ivy Peng, Ronald G. Minnich, Maya B. Gokhale: *Accelerator integration in a tile-based SoC: lessons learned with a hardware floating point compression engine.* SC 2023: 1662-1669



Lawrence Livermore National Laboratory

Maya Gokhale LLNL



Recode Engine Accelerator Unstructured Data Processor



EXAMPLE

SpMV with compressed CSR Block

Data movement is expensive both in energy and bandwidth. Recode Engine is a high-performance accelerator that can be used to compress & decompress in an energy-efficient way.



In our case, Recode-Engine decompresses Sparse Matrices for the host to compute.



Science Chiplets Opportunity (matches hand-tuned H100 performance w/12nm tech)



PCle Gen 5 · x16 Card · 64 GB HBM memory 1,6 TB/s Memory Bandwitdh · 64 GB/s IO 180 Watt Power Consumption





Benchmark Results

RTM TTI Benchmark

SYSTEM	POWER CONSUMPTION	MSTENCILS/S	
NVIDIA H100 SXM	700 WATT	21.700	
STX CARD	180 WATT	23.400	

3 0

Offload model is an unproductive way to use chiplets-base accel

(redesign for static dataflow and deep flow-through pipelines?)



Codelet Program Execution Model for Chiplets? We have to execute fast enough to run full codes

Guang Gao Jose Diaz



Design of a Codelet-based PXM with emphasis on memory orchestration

- Task and data orchestration of extremely heterogenous chipset-based systems.
- Orchestration of data movement and transformations (e.g., recoding) between tiles.



Hardware Simulation is Slow





MoSAIC: Modular System for Accelerator Integration and Communication Cross-USG Heterogeneous Integration Sabric



Lean and Mean Operates at 250MHz (1/4 real-time) Driverless inter-accelerator interaction PGAS + MsgQs for communication C++20 software stack





Customizable/Modular Open Architecture for FPGA



Open-Nic-Shell and C/C++ Software support



Non-Blocking message queue instructions similar to one-sided MPI functions: qPut(dest, data); qGet(src, data); qWait(); gPoll(); mPut(dest, data); mGet(src, data);

x_loc[i_loc] += z;

if (i == j){

else {

int j = A_j[jj];

num_qPuts++;

Open-Nic-Shell provides support for PCIE and 100Gpbs Ethernet ports.

Clock cycles per core for the asynchronous triangular solver







3

RegIO goal: host-side accessibility to the FPGA design


Integration with ChipYard/FireSim

• Chipyard is Berkeley's SOC generator platform

- NoC Generator (many topologies and bit-widths)
- Many different cores including 64-bit, OOO, and lightweight cores
- Has all of the extra bits needed to boot Linux
- Backed by 10 years of tape-outs
- FireSim
 - FPGA based simulation platform for ChipYard designs
 - GoldenGate inserts the delays to enable going from functional to cycle-accurate simulation
- MoSAIC / FireSim Integration
 - PGAS and MsgQ extensions implemented in RISC-V ROCC Interface
 - Modified Constellation NOC Generator
 - Enabled us to port MoSAIC Designs over to ChipYard and FireSim



More Efficient Chiplet Development and Integration Path



Chiplets at LBL

http://chiplets.lbl.gov LBNL/OCP Open Chiplet Economy Experience Center



Hosted by Lawrence Berkeley National Laboratory (LBNL)

Co-organized by the Open Compute Project (OCP)

Date: June 24, 2024

Time: 12:00pm to 5:00pm

Location: Berkeley National Lab, Wang Hall Bldg. 59, Room 59-3101

2024 JEDEC and OCP Standards for Chiplet Design with 3DIC Packaging Workshop



Final Thoughts: How to engage



Launching new workstream for Modular AI and HPC. Join us!



P38 Prototyping Framework



- Overall Prototyping Framework with MoSAIC integration platform takes us from Abstract Machine Model to functional model on FPGA running a full software stack
 - MoSAIC is up-and-running on FPGA platform (192 RISC-V cores per FPGA + accelerator tiles cointegrated)
 - RTL to FPGA for accelerated hardware testing & software development
 - Synthesizable RTL tape-in for FY24 (setup for tape-out to chiplets)
- GitHub Repository: <u>https://github.com/PatriGonzalez/P38_Mosaic</u>
- Tutorial: <u>https://docs.google.com/document/d/1ZjClGr6vOfTrv2l9z723LdF-E3mwk1Dglnhbyacd5SA/edit</u>



How to Engage the Hyperscalers and Advanced Packaging Community



- **Open Compute Project : Hyperscalers Interoperable Standards**
 - Open Domain Specific Architectures & Open Chiplets Economy TWG
 - Modular AI/HPC Workstream (AMD, LBL, + many others)
 - Modularity is central to OCP's existence



HETEROGENEOUS

INTEGRATION ROADMAP

SEMI APHI : Advanced Packaging & Hetro Integration

- Direct engagement with the foundries
- OSATs (Outsourced Semiconductor Assembly and Test)
- (see Samsung 2.5d pkg. service): They are an OEM, but also an OSAT

Heterogeneous Integration Roadmap: IEEE/EPS + Industry

- Everything you ever wanted to know about AdvPkg & use cases
- **CHIPS NAPMP:** National Advanced Packaging and Mfr. Program
 - Subu Ayar (formerly UCLA) and Bapi Vinnakota (formerly LBNL)





mm

BERKELEY LAB

Final Thoughts

- <u>Conventional Wisdom</u>: In the era of the "universal computer," scale was the correct answer to deliver value to HPC's scientific customers.
- <u>New Wisdom</u>: In this post-Moore/post-Exascale era, scale alone is not a viable approach to continuing to deliver value to the scientific community!
 - It isn't system scale, system effectiveness for the workload targets (e.g. refocus on strong scaling)
 - This can be delivered by specialization, and the workflows OUTSIDE of the HPC system
- <u>Conventional Wisdom</u>: Buy off the shelf microprocessors
- <u>New (old) Wisdom:</u> Follow the money (e.g. follow what hyperscale is doing)
 - <u>Join-em</u>: Buy AI machines and port everything to AI
 - <u>Beat-em (really join-em)</u>: Work together with the hyperscalers to address capital cost challenge through *modularity* and *specialization* (its not chip cost... platform cost!!!)
 - <u>HPC punches above its weight because it engages in pre-competitive R&D with industry partners</u>



Start Talking about MoSAIC Platform



MoSAIC: Modular System for Accelerator Integration and Communication Cross-USG Heterogeneous Lategration Sabric

IFC1_top-

IFC1_right-

_IFC1_left

_IFC1_left_

IFC1_bottom

Left





RISC-V ISA extensions for message driven and PGAS computing







Message Queues QPUT (destinationQID, sourceregisterID)

- QGET (destinationQID, destinationregisterID)
- QWAIT (destinationQID, statusregisterID)
- QPOLL (destinationQID, branchTarget)

Non-Blocking 3rd party memory instructions

- MPUT (Data, AddressDest)
- MGET (AddresSource, AddressDestination)



MoSAIC supports C/C++ software stack (C++20)



Applications (Histogram, Jacobi, Triangular solver...)

Compiler infrastructure (gcc++, linker scripts, helper scripts for GAS support)

> Libraries (C++, ISA extensions)

Baremetal system

```
/******
* Start Jacobi iterations *
for (int iter = 0; iter < num_iters; iter++){</pre>
  for (int i_loc = 0; i_loc < n_loc; i_loc++){</pre>
     int i = loc_to_glob_row(i_loc, tid);
     int z = r_loc[i_loc];// / A.diag[i_loc];
     x_loc[i_loc] += z;
     for (int jj = A_start[i_loc]; jj < A_start[i_loc+1]; jj++){</pre>
        int j = A_j[jj];
        int y = A_data[jj] * z;
        if (i == j){
           r_loc[i_loc] -= y;
        else {
           MsgQ_Put(j, y); // gPut()
           num_qPuts++;
```



Examples of Scalable Tiled Arrays on a U250 FPGA



multiplier, divider, and square root



Examples of mapping these arrays to an FPGAs can fit up to 192 RISC-V tiles per U250 FPGA



7 x 4 + DRAM

7 x 8 + DRAM



Examples of exploration for tri-solve only 4x slower than real-time (operating at 250MHz)



X axis: Number of parallel cores Y axis: Number of clock cycles

Total computing time



X axis: Number of parallel cores Y axis: Number of clock cycles normalized



RegIO goal: host-side accessibility to the FPGA design



Integration with ChipYard/FireSim

- Chipyard is Berkeley's SOC generator platform
 - NoC Generator (many topologies and bit-widths)
 - Many different cores including 64-bit, OOO, and lightweight cores
 - Has all of the extra bits needed to boot Linux
 - It's a heavy lift to pick it up, but its backed by 10 years of tape-outs
- FireSim
 - FPGA based simulation platform for ChipYard designs
 - GoldenGate inserts the delays to enable going from functional to cycle-accurate simulation
- MoSAIC / FireSim Integration
 - PGAS and MsgQ extensions implemented in RISC-V ROCC Interface
 - Modified Constellation NOC Generator
 - Enabled us to port MoSAIC Designs over to ChipYard and FireSim



Technology Insertion into Mainstream Platforms

AMD, Intel, Arm offer integration path for 3rd party accelerator "chiplets"

Modular AMD Chips to Embrace Custom 3rd To 'Meteor Lake' and Beyond: How **Party Chiplets**

By Francisco Pires last updated June 20, 2022 News

Supercharging learnings - and earnings - from the console space.

🚹 🛯 🚳 😰 🕞 💟 🗭 Comments (2)

When you purchase through links on our site, we may earn an affiliate commission. Here's how it works.

GPU CPU **3rd Party** 6 AIE



Intel Plans a New Era of 'Chiplet'-**Based CPUs**

At the Hot Chips 2022 conference, Intel teased its upcoming 'Meteor Lake' and 'Arrow Lake' processor families, which will use multiple tiny tiles fused together in an attempt to break free of the limits of monolithic chip design. Here's why little tiles are a big deal.







October 19, 2023

It is safe to say that ARM isn't a scrappy startup that was once the pride of the UK. The US-based IPO made the chip designer a big-game chip player, and the new capital is kickstarting some major initiatives to find more customers for its products. A new effort called Total Design aims at making it easier for companies looking to design chips in-house, an idea gaining ground with the Al boom and chip shortages.



The Importance for High Performance DFT

BERKELEY LAB



The DFT kernel for each fragment

Communication Avoiding LS3DF Formulation – Scales O(N)



$$\sum_{j,k} \left\{ F_{222} + F_{211} + F_{121} + F_{112} - F_{221} - F_{212} - F_{122} - F_{111} \right\}$$

Building Flexible Accelerators

Parameterized Hardware Generation with CHISEL

TSQR_MC

IIII Corel

HH_Core2

HH CoreN

hh_cnt:0~n-l HH_datapat

h St: br

Timing_ctr

3-Way Parallel HH-core

hh_cnt:0~n-1

HH_datapat

St: br

hh cnt:0~n-1

HH_datapa

St: br

Timing_ctr

Timing ctr

Buffer br×n

Buffer br×n

Buffer br×n



A Single HH-core



Deeply Pipelined HH-Core



Or Generate any level of concurrency or pipelining



Building Flexible Accelerators

Parameterized Hardware Generation with CHISEL

HH_Core hh_cnt:0~n-1
Buffer br×n till_datapath St: br
Timing_ctr

Table 1.3 Resource Cost Analysis

A Single HH-core

	FP Operators			Register file (Words)		
Streaming Width	Iterative	Pipeline	Parallel (16 cores)	Iterative	Pipeline	Parallel (16 cores)
16	68	432	1,088	192	432	3,072
32	132	1,632	2,112	768	3672	12,288
64	260	6,336	4,160	3,072	28,336	49,152
128	516	24,960	8,256	12,288	222,560	196,608
256	1,028	99,072	16,448	49,152	1,764,032	786,432
S V 	treaming Vidth 16 32 64 128 256	treaming VidthIterative166832132642601285162561,028	treaming VidthIterativePipeline1668432321321,632642606,33612851624,9602561,02899,072	treaming VidthIterativePipelineParallel (16 cores)16684321,088321321,6322,112642606,3364,16012851624,9608,2562561,02899,07216,448	treaming Vidth Iterative Pipeline Parallel (16 cores) Iterative 16 68 432 1,088 192 32 132 1,632 2,112 768 64 260 6,336 4,160 3,072 128 516 24,960 8,256 12,288 256 1,028 99,072 16,448 49,152	treaming Vidth Iterative Pipeline Parallel (16 cores) Iterative Pipeline 16 68 432 1,088 192 432 32 132 1,632 2,112 768 3672 64 260 6,336 4,160 3,072 28,336 128 516 24,960 8,256 12,288 222,560 256 1,028 99,072 16,448 49,152 1,764,032



Generate any level of concurrency or pipelining



Cross Domain Optimizations Using a Uniform Tensor Notation



BERKELEY LAB

Materials Science Accelerator

Preliminary Performance for Materials Science H Ψ



Franz Francetti (CMU/FFTx) BERKELEY LAB

Customizable/Modular Open Architecture for HPDA



MoSAIC: Modular System for Accelerator Integration and Communication Cross-USG Heterogeneous Lategration Sabric





Open-Nic-Shell and C/C++ Software support



Non-Blocking message queue instructions similar to one-sided MPI functions: qPut(dest, data); qGet(src, data); qWait(); qPoll(); mPut(dest, data); mGet(src, data);

Open-Nic-Shell provides support for PCIE and 100Gpbs Ethernet ports.

Clock cycles per core for the asynchronous triangular solver





Limits of Current Practice

• Typical MPW Prototype Chip Challenges

- Big chips have huge challenges
 - Complex clock trees
 - Power distribution and power sags
 - Yield / Known good die
 - Cost per area
- MPW typically results in small chips
 - Higher yield (known good die)
 - Low bandwidth and low performance due to limited shoreline BW
 - Limited area available for peripheral support functions
- Typically require FPGA support platform for memory control & PCIe
- There is a better way with chiplets

Photonics Chiplet Proto (Columbia)



Prototype





Chiplets for Scalable MPW Prototypes





Chiplets Demystification Tutorial



Industry: Heterogeneous Integration Roadmap



HETEROGENEOUS INTEGRATION ROADMAP

2019 Edition

http://eps.ieee.org/hir

HPC and Mega-datacenters is 2nd chapter

Die + Heterogeneous

System in Package (SiP)





All future applications will be further transformed through the power of AI, VR, and AR.





What is a Chiplet - Disaggregation

Chip: Homogenous logic, the one die in a package





Chiplet: Heterogeneous logic, Many die in a package









What is a Chiplet?



Why?: Chiplets Can Lower Manufacturing and Design Costs



- **1.** Scale: Build products larger than with single die
- 2. Modularity: Partition system into recognizable blocks, buy externally
- 3. Reuse: Create variants by changing a small part of a design, faster TTM
- 4. Heterogeneous integration: Die from multiple nodes, optimized for function
- 5. Customize product by customizing chiplet

Link and Transaction Layer for BoW



Not a Free Lunch: Chiplets Make Workflow More Complex


How do chiplets enable domain specialization?

Lower cost barriers to co-integrating specialization

From DARPA CHIPS



See the multi-agency chiplets workshop at https://sites.google.com/lbl.gov/chiplets-workshop-2023/home

CHIPS modularity targets the enabling of a wide range of custom solutions

The Road to an Open Chiplet Marketplace



From A Common Language to Hardware Design



"Towards a Flexible Hardware Implementation for Mixed Radix Fourier Transforms", Mario Vega, Xiaokun Yang, John Shalf, HPEC 2023 Work done by Tan Nguyen @ LBL



Lessons from Attack of the Killer Micros



Attack of the killer micros John Markoff, May 6, 1991





It is not good enough anymore to understand the technology

Now we must also understand the market context

