

# Is it the end of Moore's Law?

"I predict Moore's Law will never end. That way I will only be wrong once!"

Alan Kay: Communications of the ACM 1989



## The fundamental problem with Wires (and data movement):

Moore's Law undermined by data movement (when smaller is not better)

#### **Energy Efficiency of copper wire:**

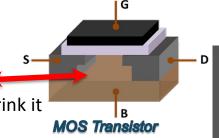
- Power = Frequency\* Length / cross-section-area



Wire efficiency does not improve as feature size shrinks

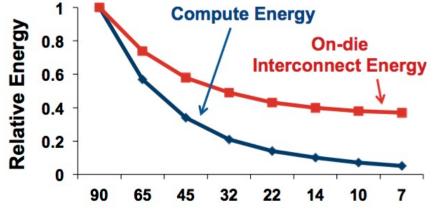
#### **Energy Efficiency of a Transistor:**

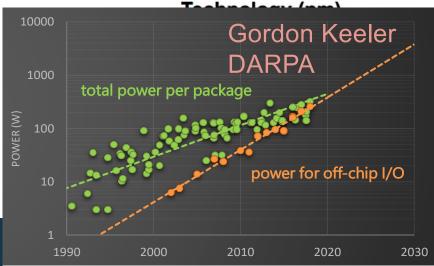
- Power = V<sup>2</sup> \* frequency \* Capacitance
- Capacitance ~= Area of Transistor
- Transistor efficiency improves as you shrink it



Net result is that moving data on wires is starting to cost more energy than computing on said data

(see also Silicon Photonics)

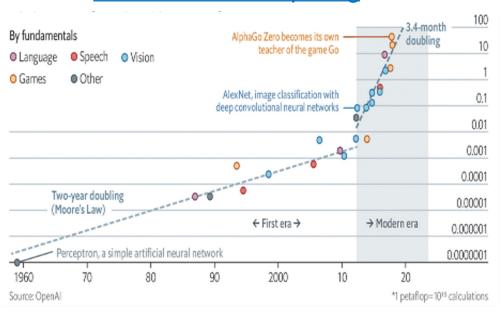




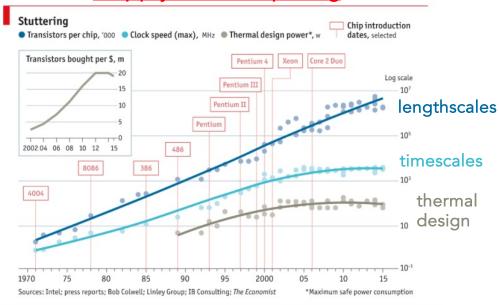


# **Explosion of Computing Demand:** Driving Need for Hyper-exponential Improvement in Performance, Energy Efficiency and Integration

#### **Demand for Computing**



#### **Supply for Computing**





## **NVIDIAnomics**



Nvidia Economics: Make \$5-\$7 for Every \$1 Spent on GPUs

By Agam Shah

- NVIDIA A100: 250W TDP
- NVIDIA H100 SXM has a 700W TDP
- Next Generation B100 is projected to consume 1400W TDP
  - Street price \$20k-\$30k
    - Prices lower @ volume
- And still supply cannot keep up with demand



# ENERGY EFFICIENCY - UNSUSTAINABLE

Al Training Energy footprint on par with entire industrial nations





Meta Al cluster 53-561 TWh\*





Ireland 31 TWh\*

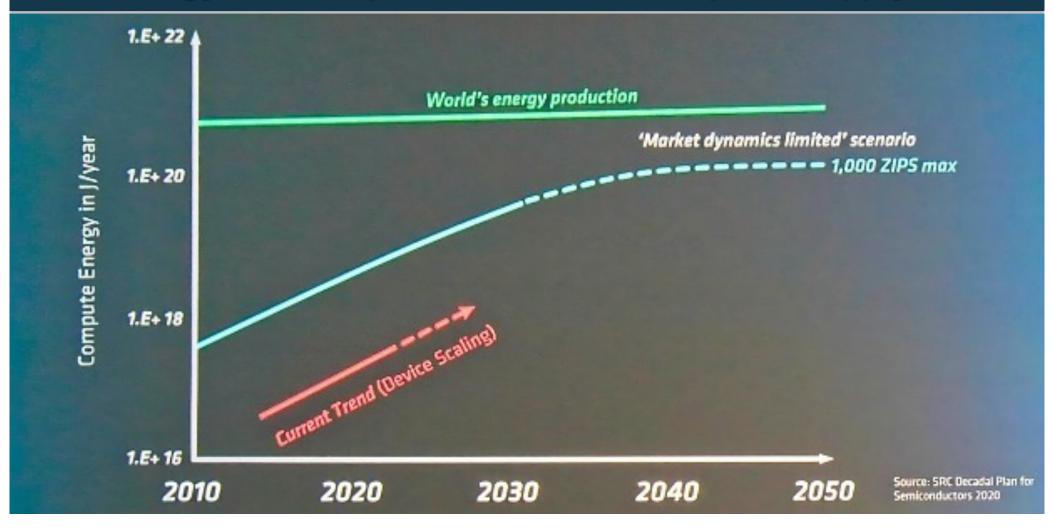


Performance Per Watt: A Critical Metric in the Age of Al

is thing. Tarytus University, Feb 2016 https://www.nemintle.com/bio/2002/05/2/1014 harming-ear-plans



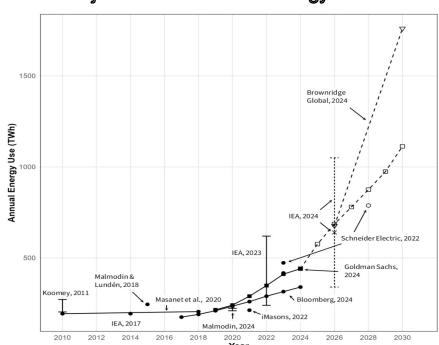
# **Al Energy Consumption On Pace to Surpass Supply**



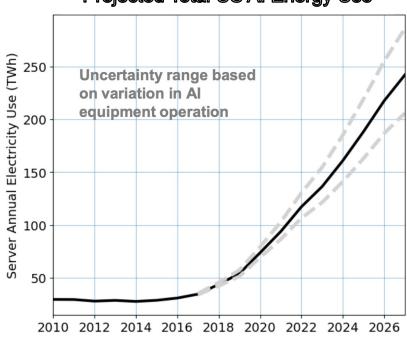
#### Impact of Al Operational Assumptions on Total Server Estimates

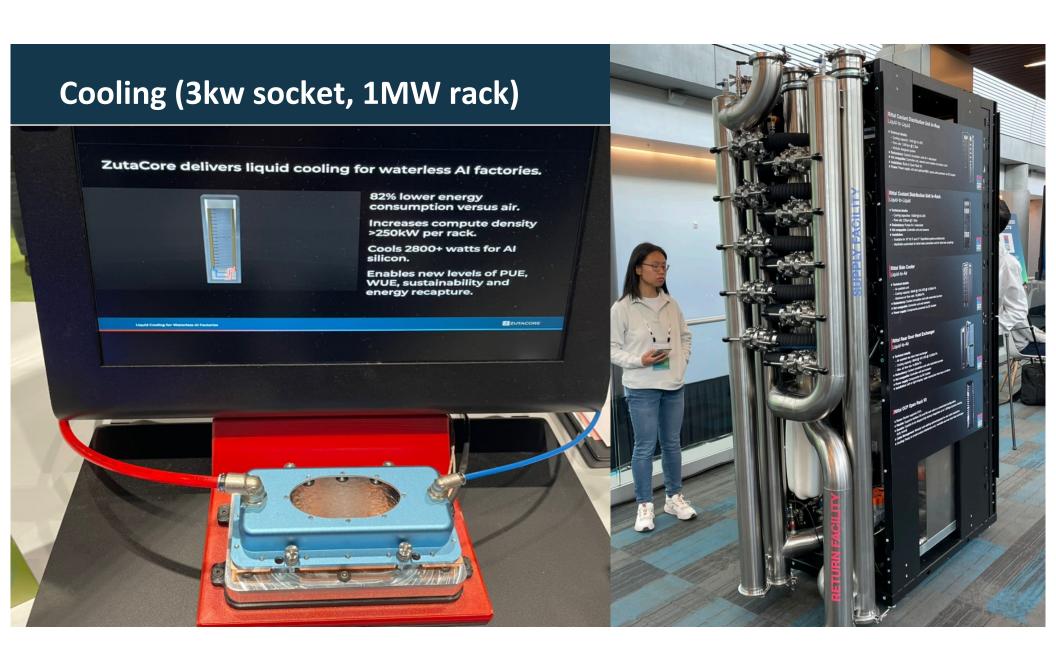
https://datacenters.lbl.gov/sites/default/files/EnergyUsageWebinar12062016.pdf

#### Projected Worldwide Al Energy Use



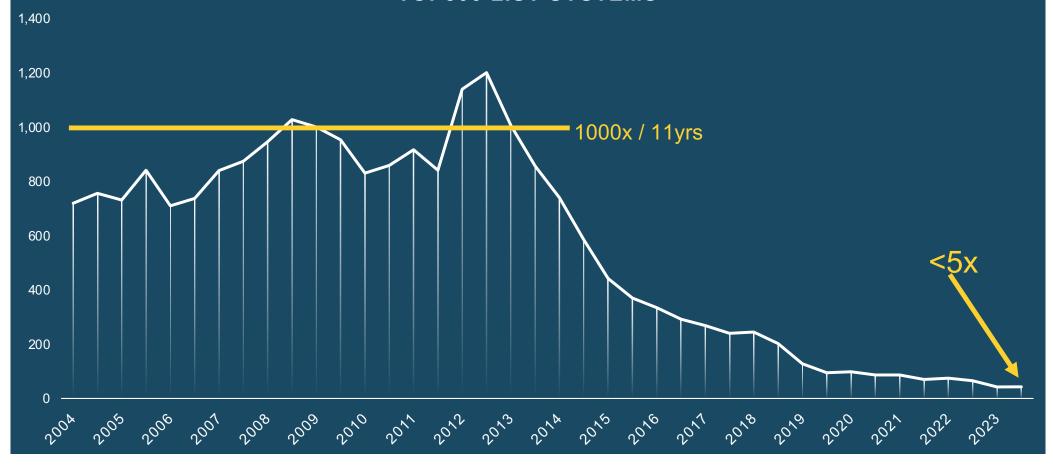
#### Projected Total US AI Energy Use



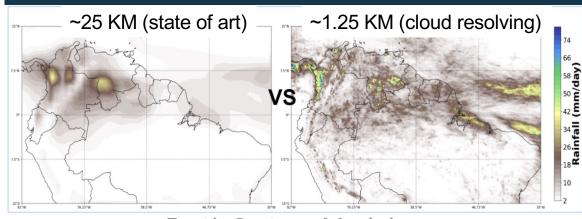


# This is HPCs future if we continue business as usual! ... and scale alone is just power and capital cost...

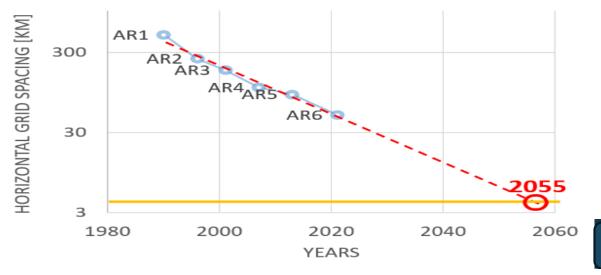
AVERAGE PERFORMANCE IMPROVEMENT PER 11 YEARS FOR SUM OF TOP500 LIST SYSTEMS



# **Example:** Kilometer Scale Climate Modeling



Earth System Models



Landmark 3.5KM Simulation on Frontier (exascale) achieved 1.5 simulated years per day performance.

At that rate, for an ensemble calculation it would take ~20 years of dedicated computing to answer important policy questions necessary to achieve 2055 goals.

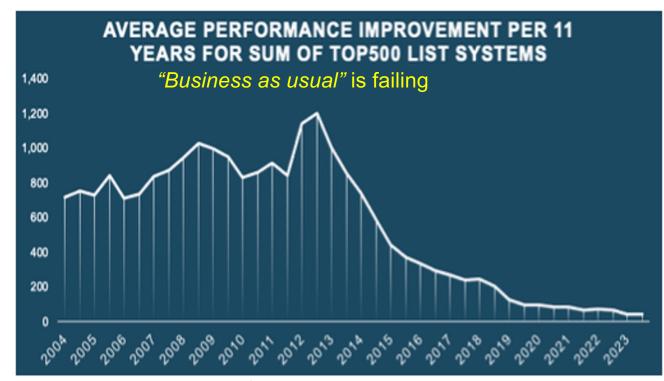
(and that is just for one policy scenario!!!)

Even if we wait for HPC performance improvements, projected 1km modeling goal will be achievable in 2055 at the current rate of progress

This is NOT an acceptable future when there are important scientific imperatives that have **global societal consequences...** 

https://climatecomputer.ncsa.illinois.edu/

# 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."

This is a call for co-design at a much deeper level than we are currently realizing

Addressing cost is crucial: Establish collaborations with industry and the Hyperscale Architects from get-go.



ASCR Competitive Portfolios in Computer Science: Algorithm Driven Codesign

# ASCR Competitive Portfolios: Algorithm-Driven Codesign of Specialized Architectures for Energy-Efficient HPC







George Michelogiannakis



Doru Popovici



Luisa Patricia Gonzalez



Kazutomo Yoshii *ANL* 



Sophia Shao UC Berkeley



Xiaokun Yang
U. Houston-Clearlake



Ann Almgren co-lead



Daniel Martin



Andrew Myers



Weiqun Zhang



Damian Rouson



#### Building on team's expertise in:

- Computer Architecture including hardware design, and hardware generators (CHISEL and GEMMINI accelerator generator), with successful chip tape-outs including compute-in-sensor, mixed-signal and HPC-class prototypes
- Computer Languages and Compiler Systems including ML-assisted code generation and Verified Lifting
- Applied Mathematics including algorithms coupling models suitable for different scales; Adaptive Mesh and Algorithm
   Refinement (AMAR) paradigm, PDEs on structured grids, spectral methods (FFTs), particle methods, and scalable solvers



## **Specialization:**

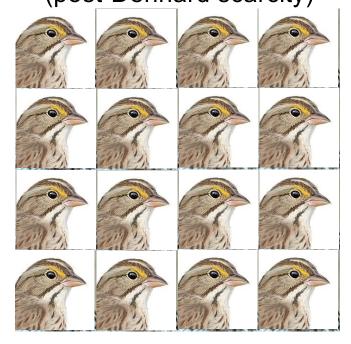
## Natures way of Extracting More Performance in Resource Limited Environment

#### **Powerful General Purpose**

# **Many Lighter Weight** (post-Dennard scarcity)

# Many Different Specialized (Post-Moore Scarcity)







Xeon, Power

KNL AMD, Cavium/Marvell, GPU

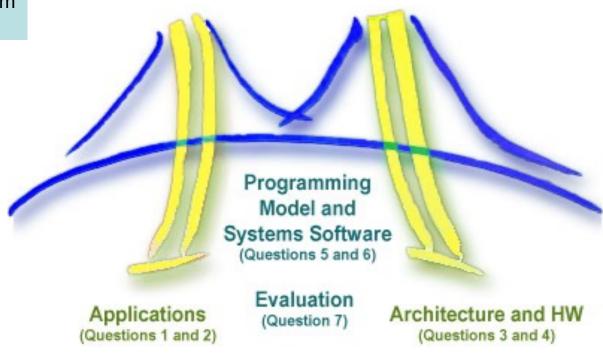
Apple, Google, Amazon, Microsoft Azure



# Phil Colella's 7 Dwarfs of Scientific Computing High-end simulation in the physical sciences = 7 numerical methods:

Exploit the mathematical structure of the problem design principle for "analogous" computing

- Structured Grids
- **Unstructured Grids**
- **Fast Fourier Transform** 3
- Dense Linear Algebra
- Sparse Linear Algebra
- **Particles**
- **Monte Carlo**



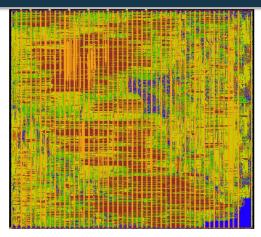


Slide from "Defining Software Requirements for Scientific Computing", Phillip Colella, 2004 Also in "The Landscape of Parallel Computing Architecture: A view from Berkeley" 2008

http://www2.eecs.berkeley.edu/Pubs/TechRpts/2006/EECS-2006-183.pdf

## **Example of Mixed-Radix MultiDimensional FFT Accelerator**

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

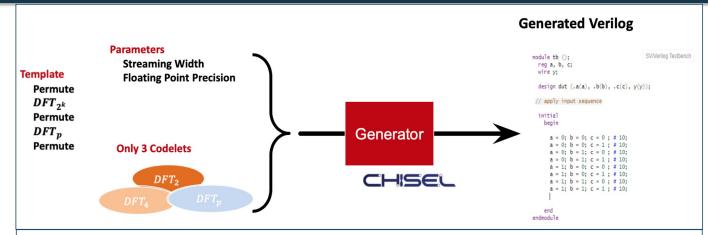


FFT96 Accelerator Die: 16nm TSMC

Area eff.: 4.18 TF/mm<sup>2</sup> Energy eff.: 4.8 TF/W

NVIDIA H100 in 4nm TSMC (10x denser than 16nm)

Area eff: 0.08 TF/mm<sup>2</sup> Energy eff: 0.0957 TF/W

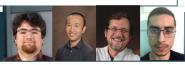


**Accomplishment:** Demonstrated general FFT Accelerator tile generator in 16nm TSMC that outperforms NVIDIA H100 by 50x in raw performance/area and 50x in energy efficiency/flop

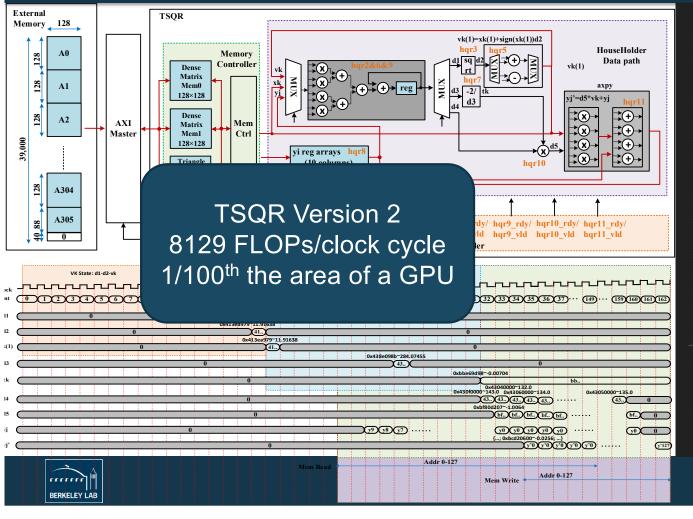
-Required only 3 codelet primitives

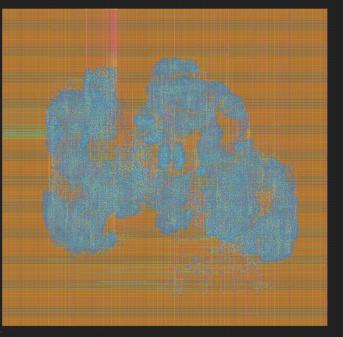
**Goal:** Generalize this approach to generators for accelerator hardware primitives that cover broad spectrum of algorithms (e.g., FFT, Dense/Sparse Linear Algebra, particles and PDEs)





# Next Steps – Implement as ASIC (OpenROAD)





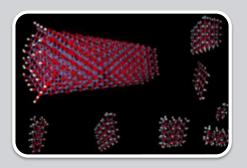
Chip: FFT9\_fpuv3
Node: freepdk45

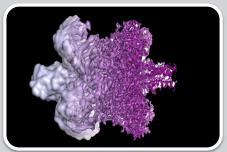
Area: 4723210.000um^2

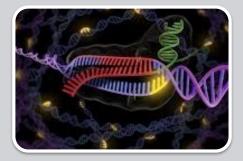
Fmax: 390.346MHz

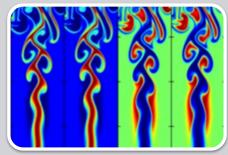
# **Architecture Specialization for Science**

(hardware designed for the algorithms) can't design effective hardware without math









#### **Materials**

Density Functional
Theory (DFT)
Use O(n) algorithm
Dominated by FFTs
(1/3 of DOE workload)

#### **Smart Sensors**

CryoEM detector 1Terabyte / sec

Compute at data source (data2info)

#### **Genomics**

String matching

Ultra-narrow datapath 2-8bit (ACTG)

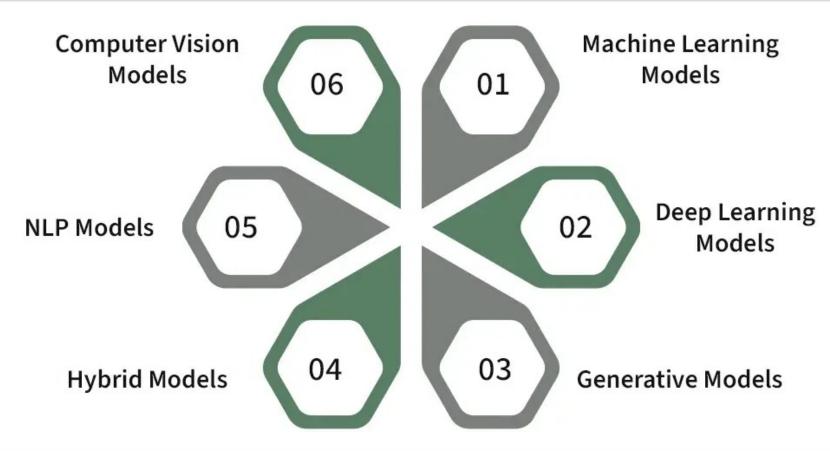
#### PDEs on Block Struct. Grids

**Extreme Strong** 

Scaling for CFL limited problems
(bigger is not better)



# There are in fact *many different kinds* of AI models with significantly different requirements





## **Neil Thompson: Economics of Post-Moore Electronics**





#### The Top

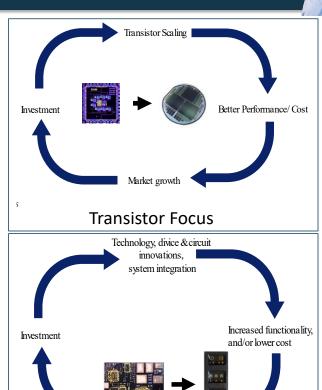
Technology	01010011 01100011 01101001 01100101 01101110 01100011 01100101 00000000	Algorithms	Hardware architecture
Opportunity	Software performance engineering	New algorithms	Hardware streamlining
Examples	Removing software bloat Tailoring software to hardware features	New problem domains New machine models	Processor simplification  Domain specialization

#### **Papers**

#### The Bottom

for example, semiconductor technology

- 1. The Economic Impact of Moore's Law
- 2. There's Plenty of Room at the Top: What will drive computer performance after Moore's Law?
- 3. The Decline of Computers as a General Purpose Technology



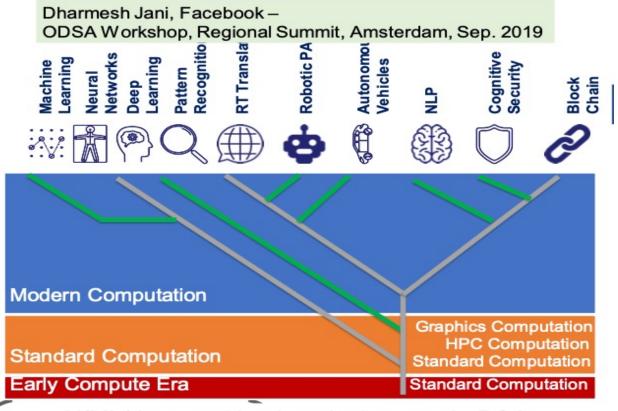
**System Focus** 



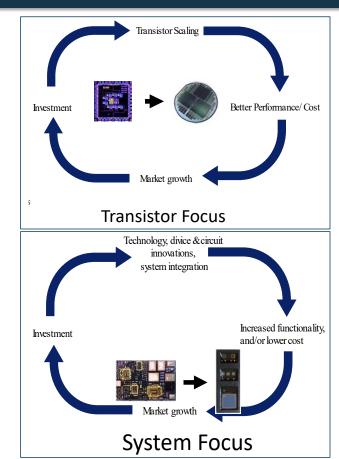


### **Domain specific Architectures driven by hyperscalers**

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



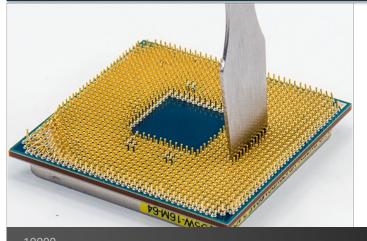
Al/ML/data workload explosion needs DSAs

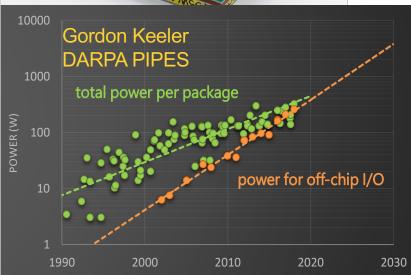




## Why Advanced Packaging & Chiplets?

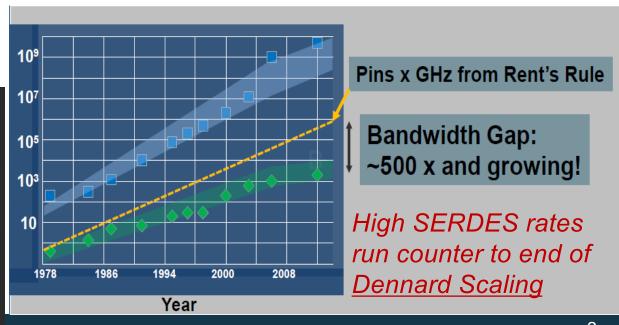
Because Package Performance is Pin Limited





#### **Rent's Rule:**

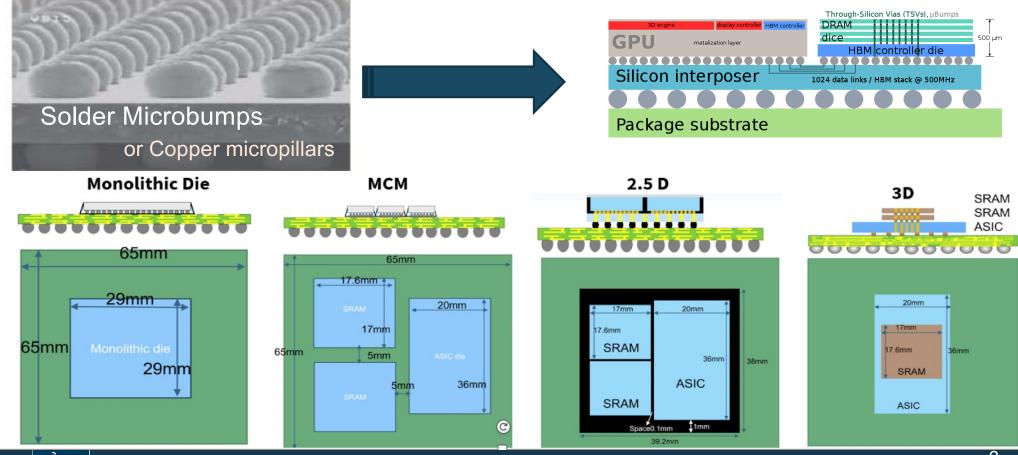
Number of pins = K x Gates<sup>a</sup> (IBM, 1960) K = 0.82, a = 0.45 for early Microprocessors



Source: J. Poulton, Nvidia

# What is a Chiplet?

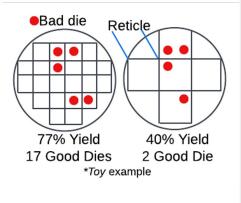
BERKELEY LAB



2 3

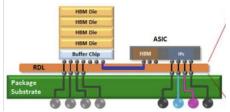
# Chiplets have many well known advantages

# Yield and cost efficiency



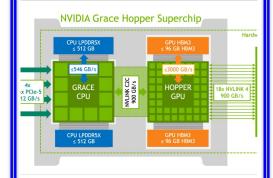
# Overcome the memory wall

HBM is 3D integration



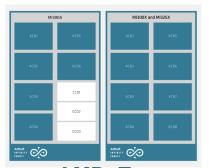
Lau. et al. Chiplet design and heterogeneous integration packaging

# Heterogeneous integration



CPU + GPU

# Scalability & Modularity/Reuse



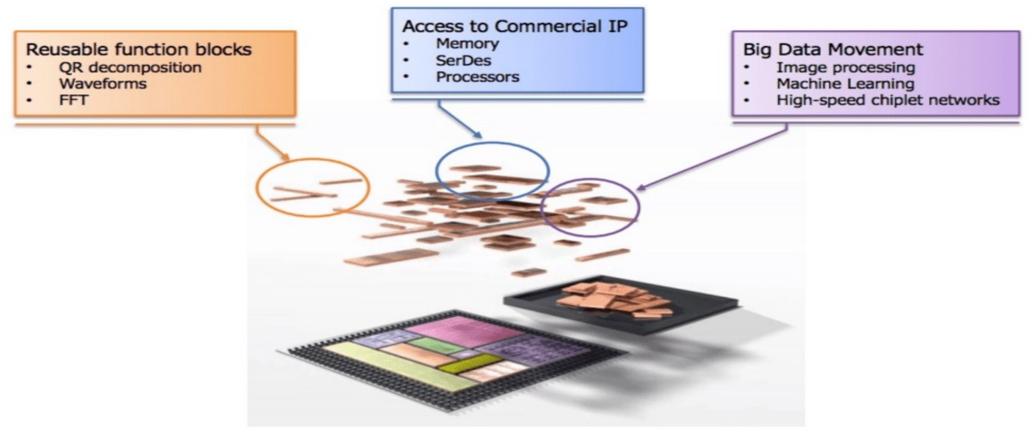
AMD. Two different products.
Same chiplets



## **How** do chiplets enable domain specialization?

Lower cost barriers to co-integrating specialization

From DARPA CHIPS

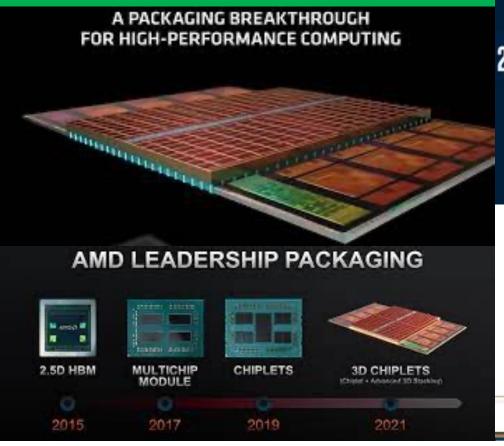


See the multi-agency chiplets workshop at https://sites.google.com/lbl.gov/chiplets-workshop-2023/home

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

## **Heterogenous Integration Driving Exascale HPC**

AMD Advanced Packaging for OLCF Frontier



#### Intel Ponte Veccio for ALCF Aurora

COMPUTE DENSITY

**Building the Foundation for Exascale Computing** 



## **Technology Insertion into Mainstream Platforms**

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

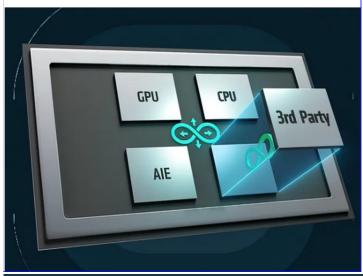
# Modular AMD Chips to Embrace Custom 3rd Party Chiplets To 'Meteor Lake' and Beyond: How

News By Francisco Pires last updated June 20, 2022

Supercharging learnings - and earnings - from the console space.



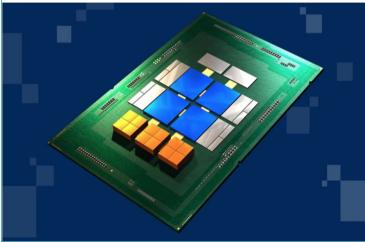
When you purchase through links on our site, we may earn an affiliate commission. <u>Here's how it</u> works.



#### To 'Meteor Lake' and Beyond: How 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.

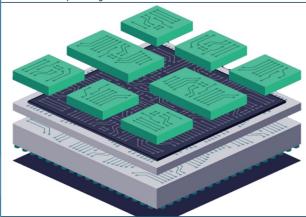






October 19, 2023

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.

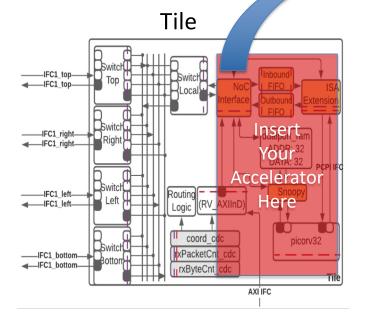




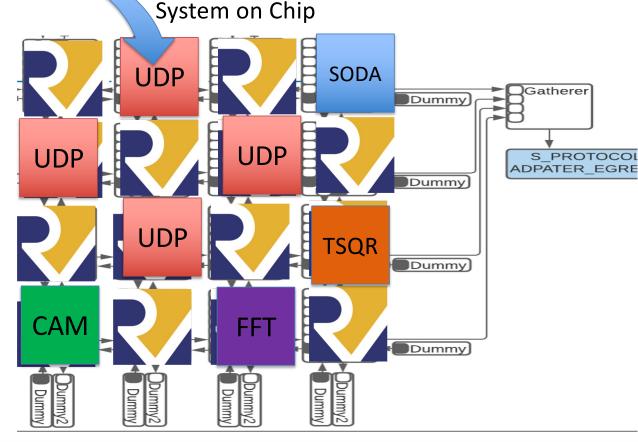
http://chiplets.lbl.gov/

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





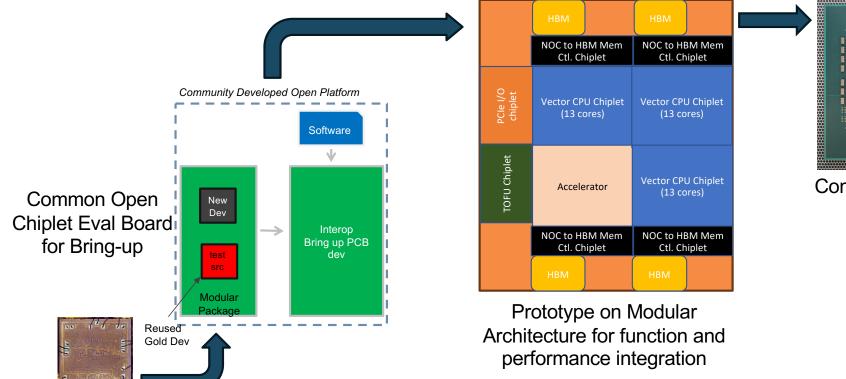
Lean and Mean
Operates at 250MHz (1/4 real-time)
Driverless inter-accelerator interaction
PGAS + MsgQs for communication
C++20 software stack





## More Efficient Chiplet Development and Integration Path

Platform for open development with path into commercial platform





Commercial product integration





# **Barriers to Industry Adoption**

# The Road to an Open Chiplet Marketplace

It's a big vision, but we are not there yet....

**Proprietary Chiplet Platform** 

> + Single vendor controlled chiplet and packaging stack

Initial Driver for Re-usable Chiplets

Semi-Custom Multi-Vendor **Chiplet Platform** 



+ PHY, Transport and Protocol standards required for interoperability

**Open Chiplet** Marketplace



- + Viable Multi-Vendor Business Model
- + PHY, Protocol and Transport
- + Mechanical, thermal & power
- Pre/post-silicon test & debug
- Software standards
- Silicon qualification, reliability and manufacturing

From Jeff DiFilippi at ARM **OCP Global Summit 2023** 

Products Markets

**TECHNOLOGY TOPICS** 

Chiplets

Today





## **Industry:** Heterogeneous Integration Roadmap



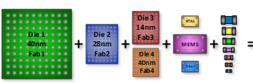


Cloud

**Data Centers** 

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

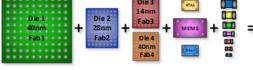
http://eps.ieee.org/hir

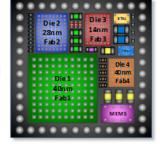


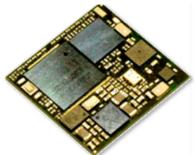
Mobile

Everywhere

HPC and Mega-datacenters is 2<sup>nd</sup> chapter







System in Package (SiP)







to

IoE



# 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





## OCP Modularity for HPC & Al Workstream needs your input!



Patricia G. (LBNL)

George M. (LBNL)

Anu R. (Microchip Inc.)

luisa.gonzalez@ocproject.net
georgios.michelogiannakis@ocproject.net
tony.gutierrez@ocproject.net
anu@microchip.com

#### How to participate:

- Subscribe to the Project Mailing List
- Add the Project Meeting Calendar
- Join the Conversation





## **Dynamic Reconfiguration (aka resource Disaggregation)**

Ultra-Performance co-packaged optics/networking for resource disaggregation



#### Diverse Node Configurations for Diverse Workload Resource Requirements

#### **Training**

- 8 connections: GPU
- 8 links to HBM (weights)
- 8 links: to NVRAM
- 1 links: to CPU (control)

GPU0 GPU3 GPU4 GPU7 GPU5 GPU5 GPU5 TOR

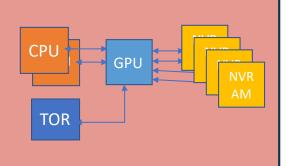
**GPU** 

#### **Data Mining**

- 6-links: HBM
- 15 links: NVRAM

(capacity)

4 links: CPU (branchy code)

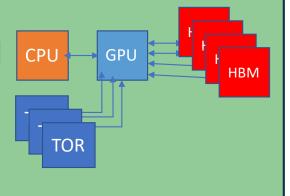


#### Inference

- 16 links to TOR (streaming data)
- 8 links HBM (weights)
- 1 link: CPU

#### **Graph Analytics**

- 16 links HBM
- 8 links TOR
- 1 Link CPU









**CPU** 



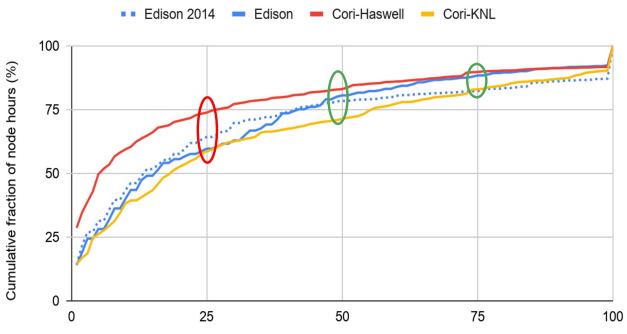
TOR





# **Need for Memory Disaggregation**

#### Memory pressure at NERSC, 2018



About 15% of NERSC workload uses more than 75% of the available memory per node.

And ~25% uses more than 50% of available memory.

But 75% of Haswell job hours (60% of KNL) use < 25% memory

Fraction of Node Memory Used (%)

Overestimate: maxrss x

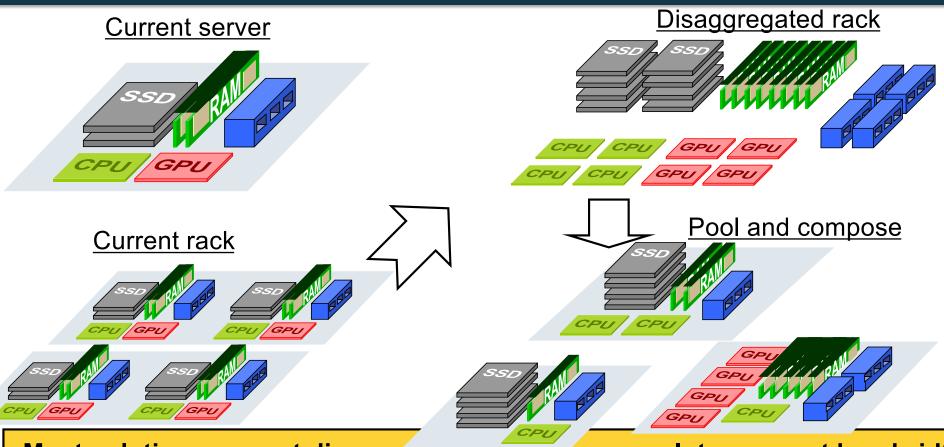
ranks\_per\_node\_

Assumes memory balance across



Brian Austin: NERSC Workload Analysis

## **Disaggregated Node/Rack Architecture**



Most solutions current disaggregation solutions use Interconnect bandwidth (1 – 10 GB/s)

## **An Al Chiplets Example in 6 Easy Steps**

Memory Disaggregation within the Package

Luisa Patricia Gonzalez Guererro (LBNL)





# Recipe for Hardware Specialization enabled by chiplets modularity in 6 steps

#### Ingredients:

- 4 Chiplets
- 1 Interposer

#### **Chiplets Marketplace**;)



CNNs
Transformers

DeepSpeech 2

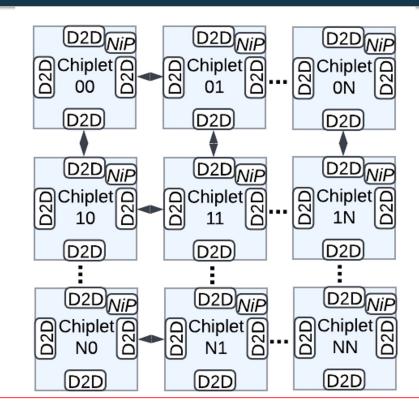
Three optimal Systems in Package (SiP) for Al



## Step 2: Scalable Interposer 2D grid with NxN Chiplets

#### Network-In-Package (NiP)

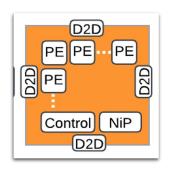
- Packet switched
- Light weight
- Minimum transmission latency (1CC)
- No Cache Coherency
  - Message passing for parallel computing
    - Hardware Message Queues

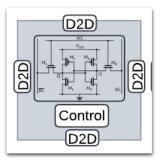


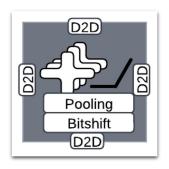
Dataflow is regular and predictable

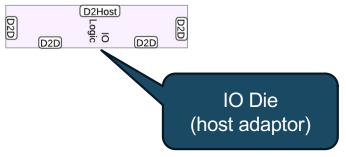


## Step 3: Let's do some math and define our chiplets parameters







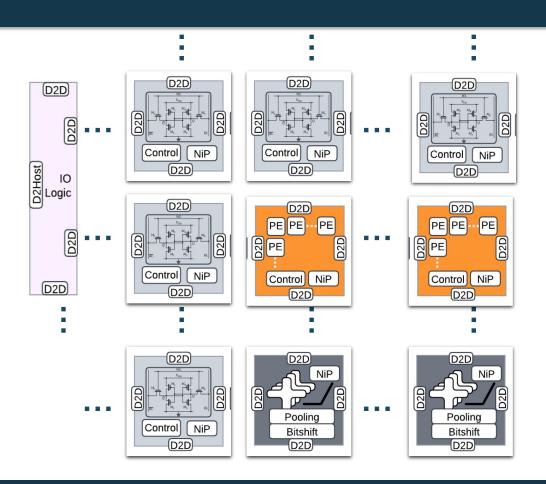


	Chiplet 4 mm x 4mm	Sub 20nm @ 1GHz	Sub 10nm @ 1GHz Sub 5nm @ 2GHz	
1	Systolic Array	64 x 64 PEs	128 x 128 PEs	
2	Scratchpad	1MB	9MB	
3	Accumulator	64 acc	128 acc	



## **Step 4: Architecture for an AI SiP**

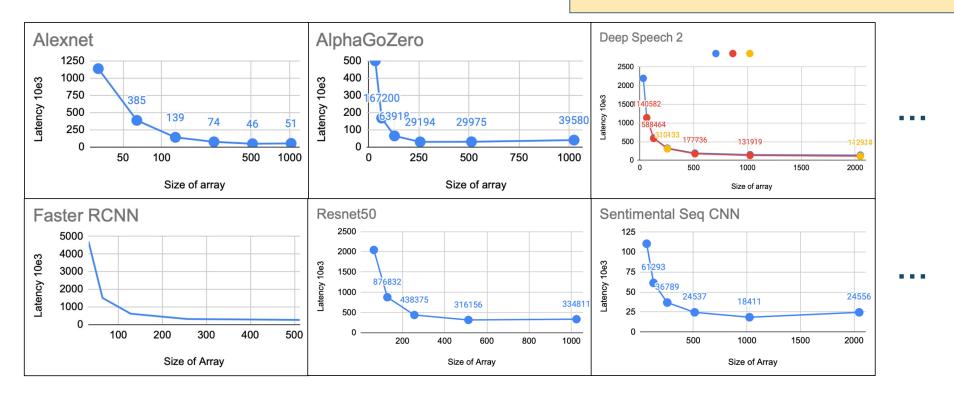
- Scratchpad chiplets
  - Left: Input/Output Buffer
  - Top: Trained parameters
- KxK systolic arrays chiplets
- Accumulator chiplets : last row
- IO Chiplet: From and to host





## **Step 5: Profiling AI algorithms**

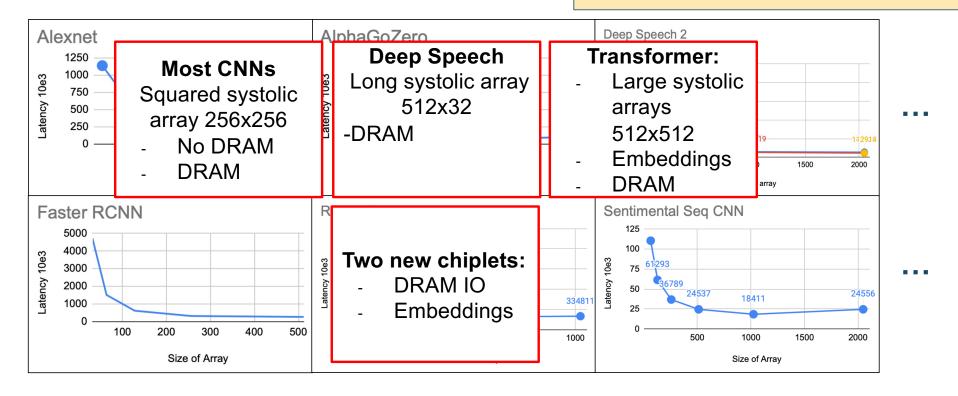
Optimal Size of systolic arrayOptimal Scratchpad capacity (MB)





### **Step 5: Profiling AI algorithms**

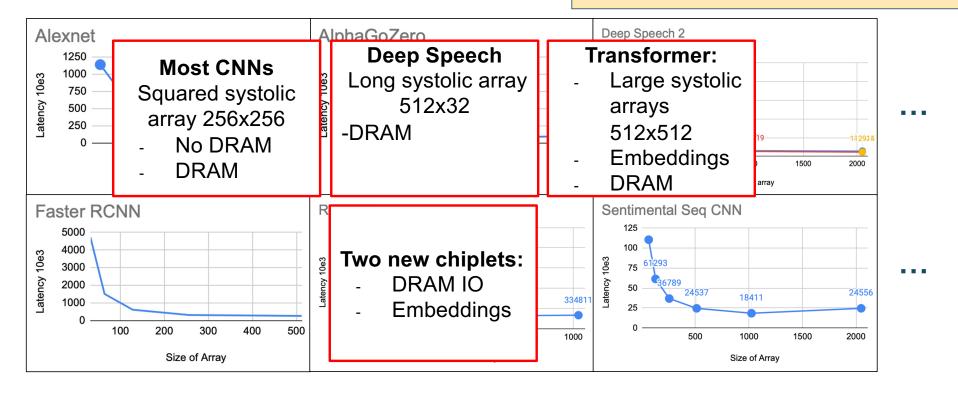
Optimal Size of systolic Array
 Optimal Scratchpad capacity (MB)





#### **Step 5: Profiling AI algorithms**

Optimal Size of systolic ArrayOptimal Scratchpad capacity (MB)





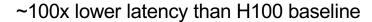
## Step 6: Three different SiPs for three dataflows

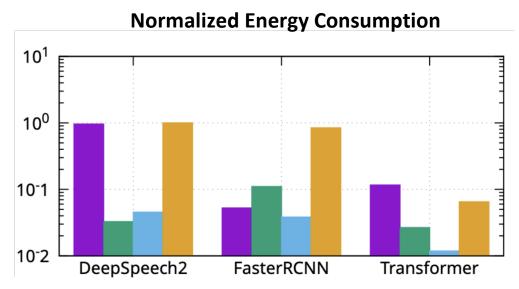
	Sub 20nm SiP			Sub 10nm/5nm SiP		
	CNN	Deep	Trans.	CNN	Deep	Trans.
		Speech			Speech	
Syst. Array PEs	256x256	512x64	512x512	256x256	512x128	512x512
Syst. Array Chiplets Num.	16	8	64	4	8	16
Scratchpad Chiplets Num. (3D)	16	27	20	10	27	12
Accumulator Chiplets Num.	4	1	8	2	1	4
Embeddings Chiplets Num.	0	0	8	0	0	4
BW [TB/s] per Chiplet	1	1	1	2/4	2/4	2/4
Data lanes per Chiplet	8192	8192	8192	16384	8192	16384
Chiplet Grid	6x6	6x6	10x10	4x4	6x6	6x6
Interposer Size* [mm x mm]	38x38	38x38	64x64	26x26	38x38	38x38



## **Evaluation: Better Latency and Better Energy consumption**

# Normalized Latency 101 100 10-1 10-2 10-3 DeepSpeech2 FasterRCNN Transformer

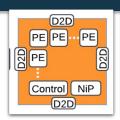


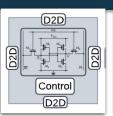


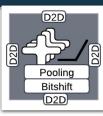
~10x lower energy consumption than H100 baseline



#### **Overspecialization: Fear Not!**









4 Chiplets + 1 interposer = 4 highly specialized SiPs for AI

+100x less latency than baseline +10x better energy consumption than baseline

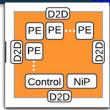
#### Final thoughts and next steps

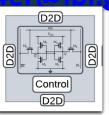
- Chiplet = Kernel, *i.e.* (GEMM)
- HPC applications:
  - Density Functional Theory
  - Molecular Dynamics
  - Climate modeling
  - Processing In Cell
  - NEGFs
  - SpVM

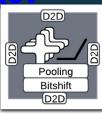


#### Thank you!

lg4er@lbl.gov









4 Chiplets + 1 interposer = 4 highly specialized SiPs for Al

100x less latency than baseline10x better energy consumption than baseline

And a lot of that was about memory specialization and memory fabrics





## **CoPackaged Optics and Photonic Disaggregation**

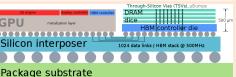
Lets take this to the net level!



## Impedance Matching to Packaging Technology



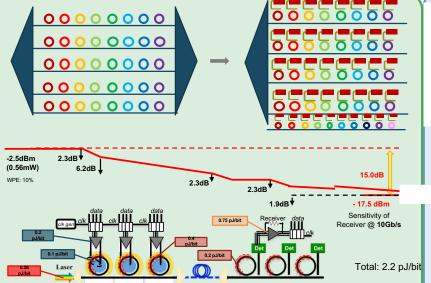




#### **In-package integration**

Solder Microbumps & Copper Pillars@~10Gbps

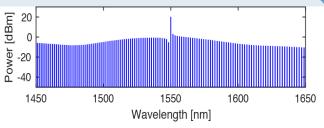
Wide and Slow!

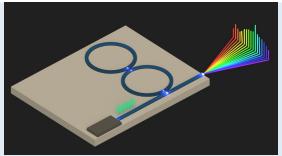


#### **DWDM Using Silicon Photonics**

Ring Resonators @ ~10-25 Gb/sec per chan Many channels to get bandwidth density

Wide and Slow!





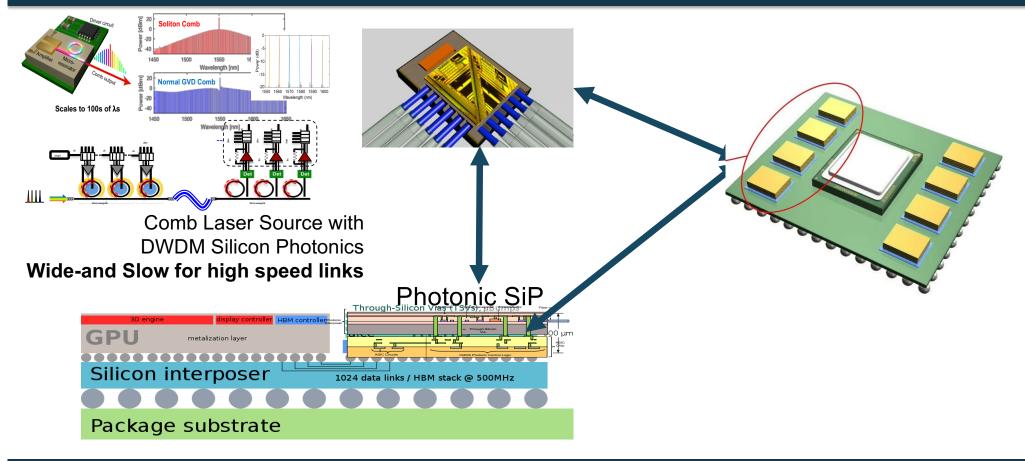
#### **Comb Laser Sources**

Single laser to efficiently generate 100s of frequencies

Wide and Slow!

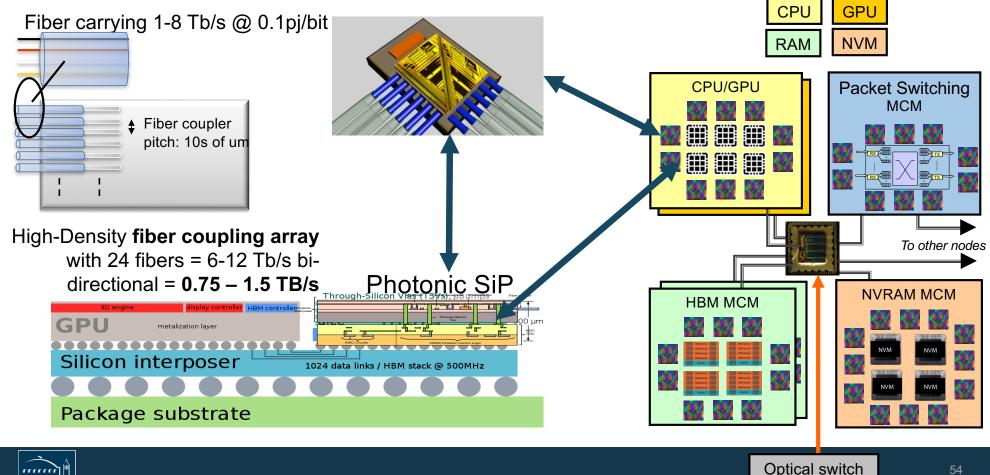


## **Photonic MCM (Multi-Chip Module)**

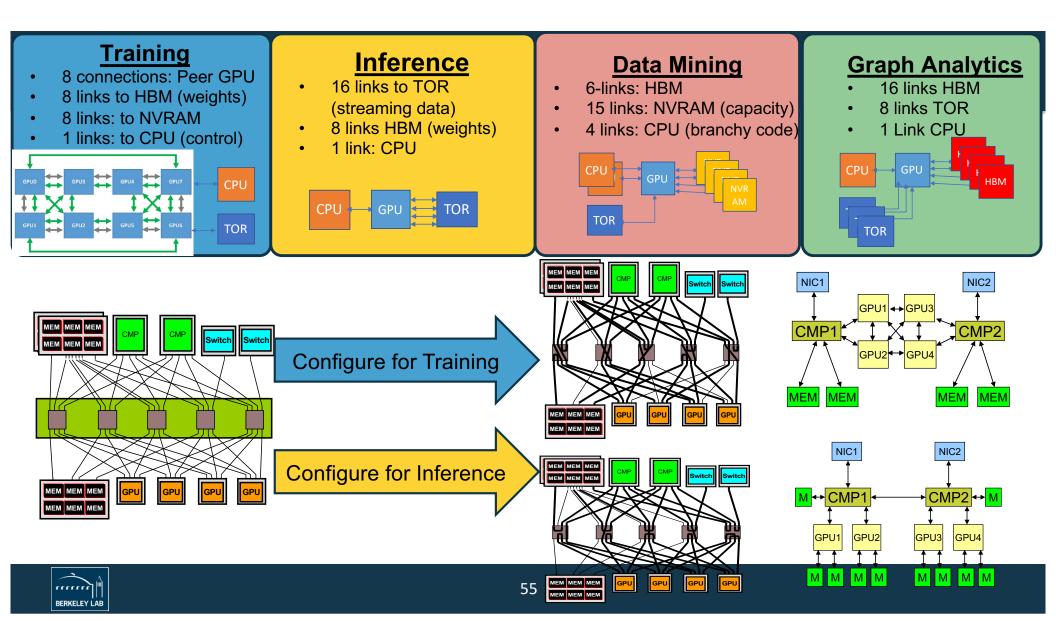




## **Photonic MCM (Multi-Chip Module)**







# Latency Sensitivity Study: Focus on Single-Hop Networks Need to minimize network diameter due to latency sensitivity

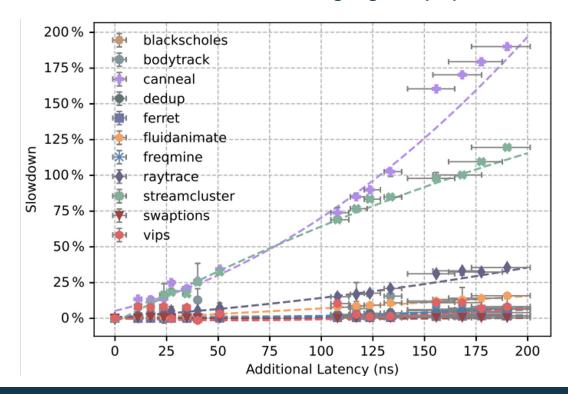
#### **PARSEC 3 on real hardware:**

- Sapphire Rapids HBM nodes
- Dual Intel Xeon Max 9462 CPUs
- Flat mode
- SNC4 clustering
- The execution core is constant
- We vary memory NUMA domain

#### Workload sensitivity to latency:

- Most show minimal effect
- Few (canneal, streamcluster) suffer significant performance degradation with latency increase

## Dots are measurements. Dashed lines are curve fitting degree 2 polynomials.



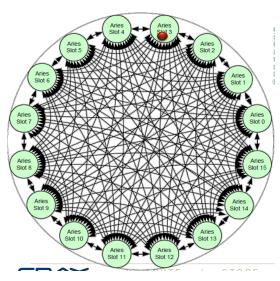


## Low-Diameter/Single-Hop Interconnects (Dragonfly or All2All)

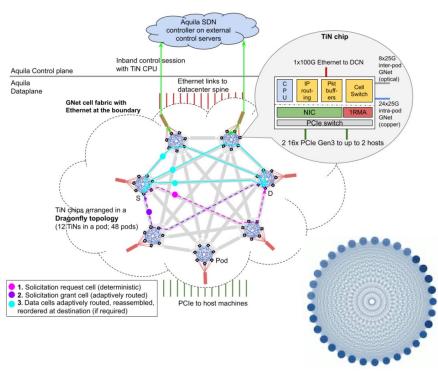
Started in HPC, but seeing traction in hyperscale and OCP

Cray Dragonfly – Aries/YARC2 2016



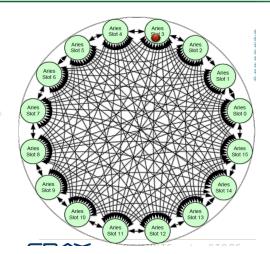


# Google Aquila 2022

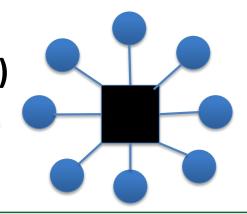


#### **Key Value Proposition for AWGR Lambda Shuffle**

- Electrical All-2-All for single-hop networks
  - $O(N^2)$  cables
  - O(N) connectors/endpoint

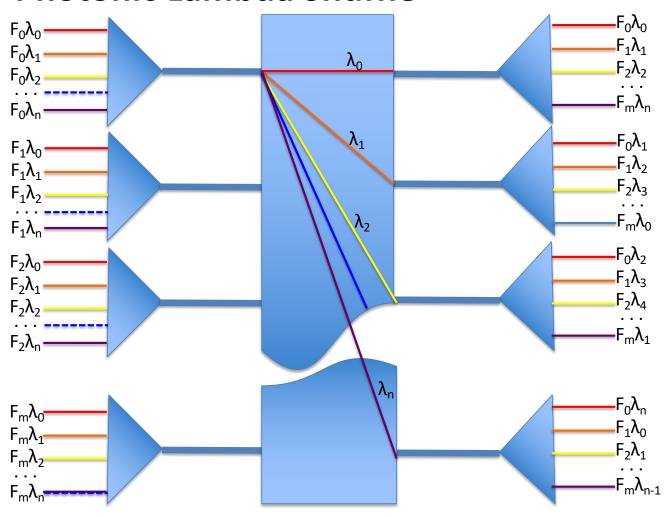


- Photonics lambda shuffle (unique to photonics)
  - O(N) cables
  - O(k) connectors/endpoint



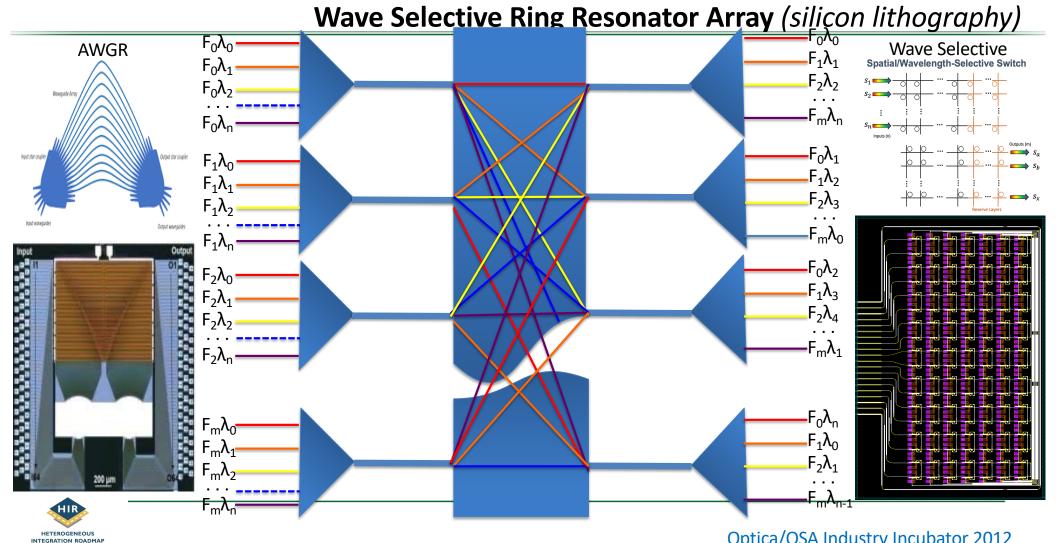


## **Photonic Lambda Shuffle**



Optica/OSA Industry Incubator 2012

#### Monolithic AWGR (stamped in plastic) or



Optica/OSA Industry Incubator 2012

## Some other stuff necessary to complete the picture

# HBM4: an HMC-like Serial Interface to HBM

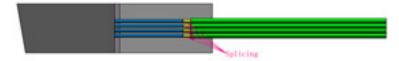
**Example Samsung Custom HBM4** 

#### Fiber Attach

Opportunity with expensive electric connectors

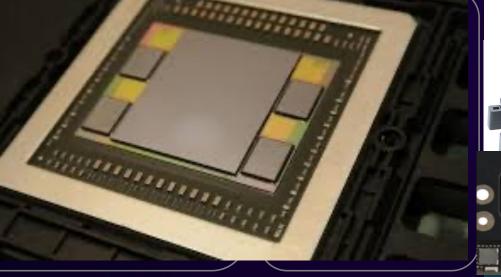
Connector (Expensive)

Fiber (cheap)

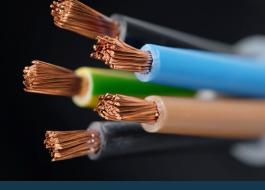


Connector (Expensive)

Copper Wire (cheap)









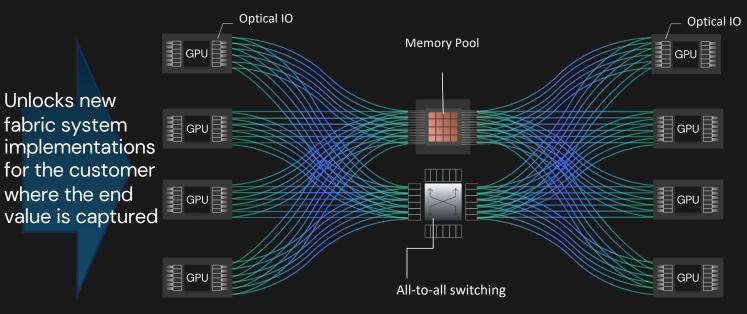


Systems solutions are <u>always</u> more powerful than any individual component Comb laser level innovation unlocks value only at Fabric System Level All solutions depend upon co-packaged optics at the endpoints

#### Comb driven Optical fabric solution



Xscape 's CombX Laser



Optical Fabric System

## **System Throughput per Dollar Spent on Memory**

Photonic disaggregation vs. conventional fixed node

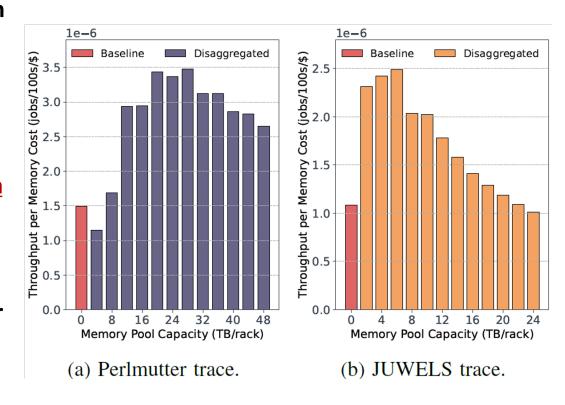
Rack-scale disaggregation with LBNL/Julich Fair Memory scheduling Policy

Assume \$4.9 per GB of memory capacity

**Baseline:** No disaggregation. 512 GB/node

Disaggregation offers 3.5x improvement in system throughput per \$

Beyond a certain memory capacity, negligible benefit in throughput but higher CapEx for memory





## Anatomy of a "Value" Metric

**Good Stuff** 

**Bad Stuff** 

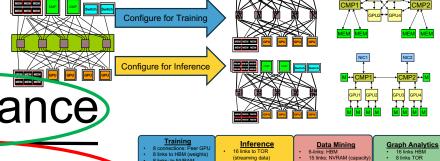


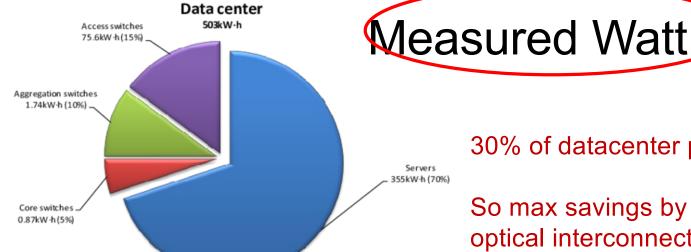
## Anatomy of a "Value" Metric

Increase performance with Disaggregation And Bandwidth Steering!

Deliver bandwidth to where it is needed. By taking it from where it isn't

Performance





30% of datacenter power goes to network

So max savings by creating infinitely efficient optical interconnect is 30%!



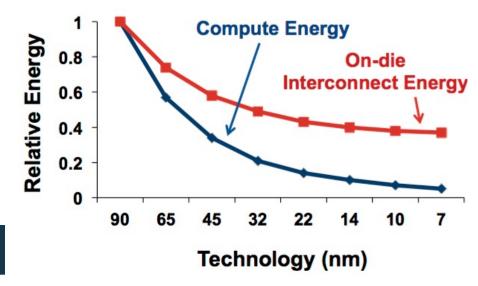
Exploit the unique properties of photonics to improve both numerator and denominator!

Some final shots and open questions for you to ponder



#### Question: If data movement is where all the energy is going?

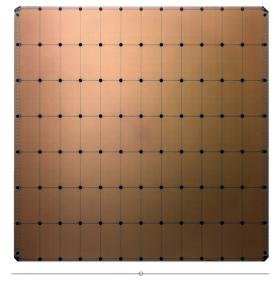
- Why do we still teach students that the order of complexity of algorithms as O(flops)?
- Can't we teach order of complexity for data movement?
  - This has been discussed in SIAM meetings for a couple of decades
  - Even started a PADAL workshop series with Thomas and Didem on this question
- What would a solution look like?





## Compute in Memory? (the elusiveness of NonVon Architecture?

- Everyone talks about Compute In Memory as if it is easy.
- The very first concept Von Neumann had for a computer had no boundary between compute in memory
  - He got over it
  - Why do we blame the poor guy for finding a solution to what is clearly a very difficult problem!
- Compute in memory is often proposed to solve this Von Neumann Bottleneck
  - Yet we have a commercially available compute in memory platform
  - Do you feel all of the problems have been solved?
  - Is it easy now?
- Also, the fundamentals of our memory is broken
  - DRAM is densest, but neither cell cycle time nor density has not improved at historical rates since hitting 10nm.





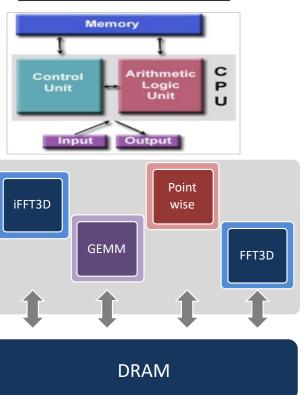




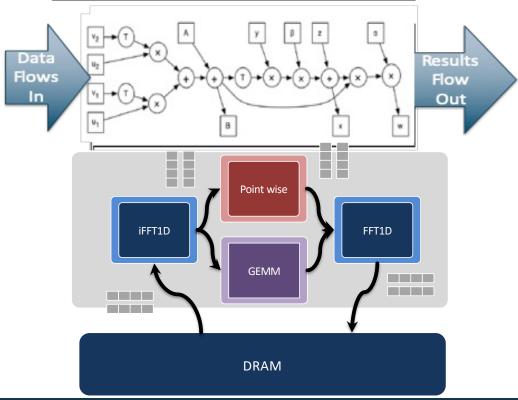
## Offload model is an unproductive way to use hetero-acceleration

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

#### Von Neumann CPU

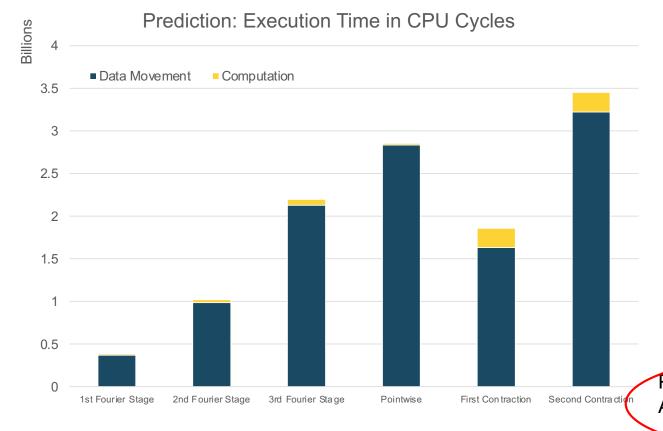


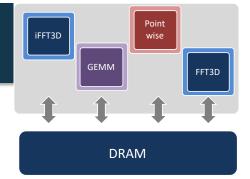
#### Dataflow (FPGA, GraphCore etc.)





## Modeling the Baseline Computation Limitations of the Offload Model for Accelerators





#### **Assumptions:**

#### A. Hardware

- CPU with 8 cores
- CPU frequency is 3.9 GHz
- DRAM bandwidth 22 GB/s

#### A. Problem Size

- Sphere diameter 64
- FFT size 128
- Number of bands 256
- Number of atoms 256

Prediction: 15.4 billion cycles

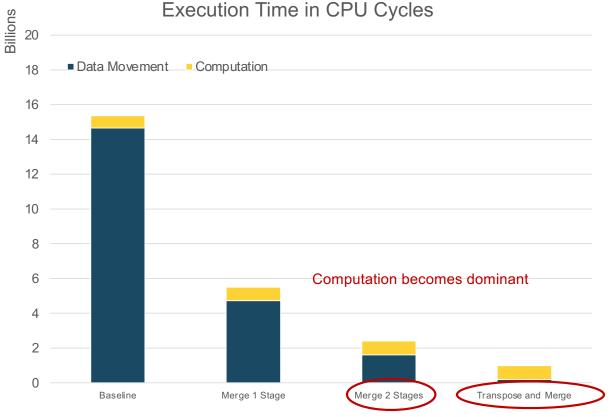
Actual execution: 17.1 billion cycles

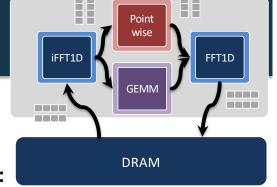
Accuracy - 90%



#### **Peer Accelerator Model (merged Kernels)**

#### Convert Problem to Compute Bound





#### **Assumptions:**

#### A. Hardware

- CPU with 8 cores
- CPU frequency is 3.9 GHz
- DRAM bandwidth 22 GB/s

#### A. Problem Size

- Sphere diameter 64
- FFT size 128
- Number of bands 256
- Number of atoms 256

Merging different stages

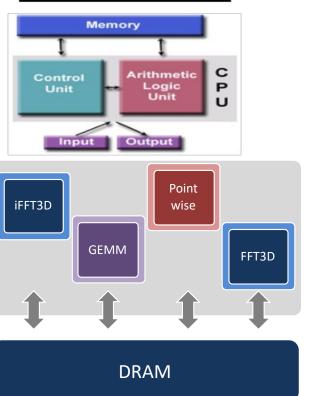
The Transpose and Merge version requires a lot of on chip memory



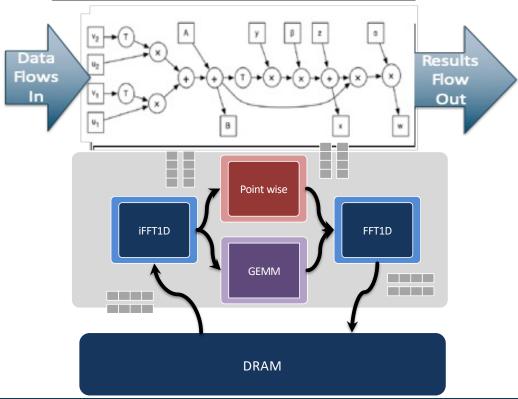
## Offload model is an unproductive way to exploit hetero-accel

(For p-models, this also breaks encapsulation in subroutines as an abstraction)

#### **Von Neumann CPU**



#### **Dataflow (FPGA, GraphCore etc.)**





This is an NP-hard graph embedding problem

#### **Conclusions**

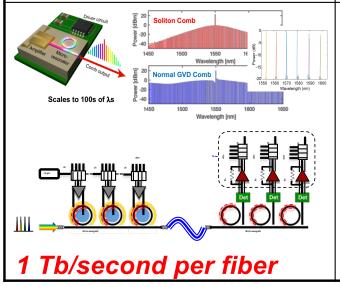
- Future of energy efficient HPC and AI lies in specialization & Systems
  - Modularity through advanced packaging
  - Codesign that spans Algorithms, Software, and Hardware
- Systems Performance (not just device performance) is the future
  - Components are bottom of the value chain
  - Systems enabled by the unique properties of photonic devices are the <u>top</u> of the value chain
  - Memory fabrics are the essential glue that enables us to delver such systems
- Requires a modular approach to deliver systems
  - Chiplets with baked-in memory fabrics are an opportunity for scalable modularity
  - OCP Workstreams to negotiate standards to promote that modularity

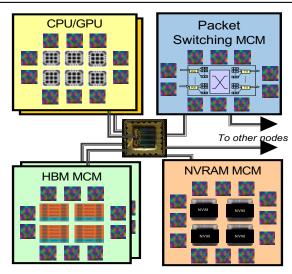


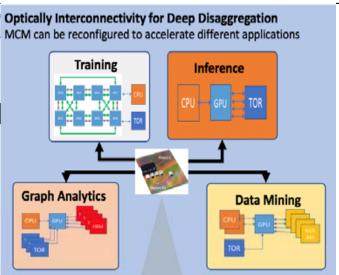
## **PINE: Photonic Integrated Networked Energy Efficient Datacenters**

Addressing the datacenter energy challenge!

1) Energy-bandwidth optimized optical links 2) Embedded silicon photonics into OC-MCMs 3) Bandwidth steering for **Custom Node Connectivity** 







**ENLITENED** 











Glick





Lipson



սիսիս

CISCO



**NVIDIA** 











#### OCP Modularity for HPC & Al Workstream needs your input!



Patricia G. (LBNL)

George M. (LBNL)

Anu R. (Microchip)

luisa.gonzalez@ocproject.net
georgios.michelogiannakis@ocproject.net
tony.gutierrez@ocproject.net
anu@microchip.com

#### How to participate:

- Subscribe to the Project Mailing List
- Add the Project Meeting Calendar
- Join the Conversation





## end

