

Exploring Emerging Technologies in the Extreme Scale HPC Co-Design Space

Jeffrey S. Vetter

11th International Conference on Parallel Processing and Applied Mathematics (PPAM)

Krakow

8 Sep 2015

OAK RIDGE NATIONAL LABORATORY
MANAGED BY UT-BATTELLE FOR THE DEPARTMENT OF ENERGY

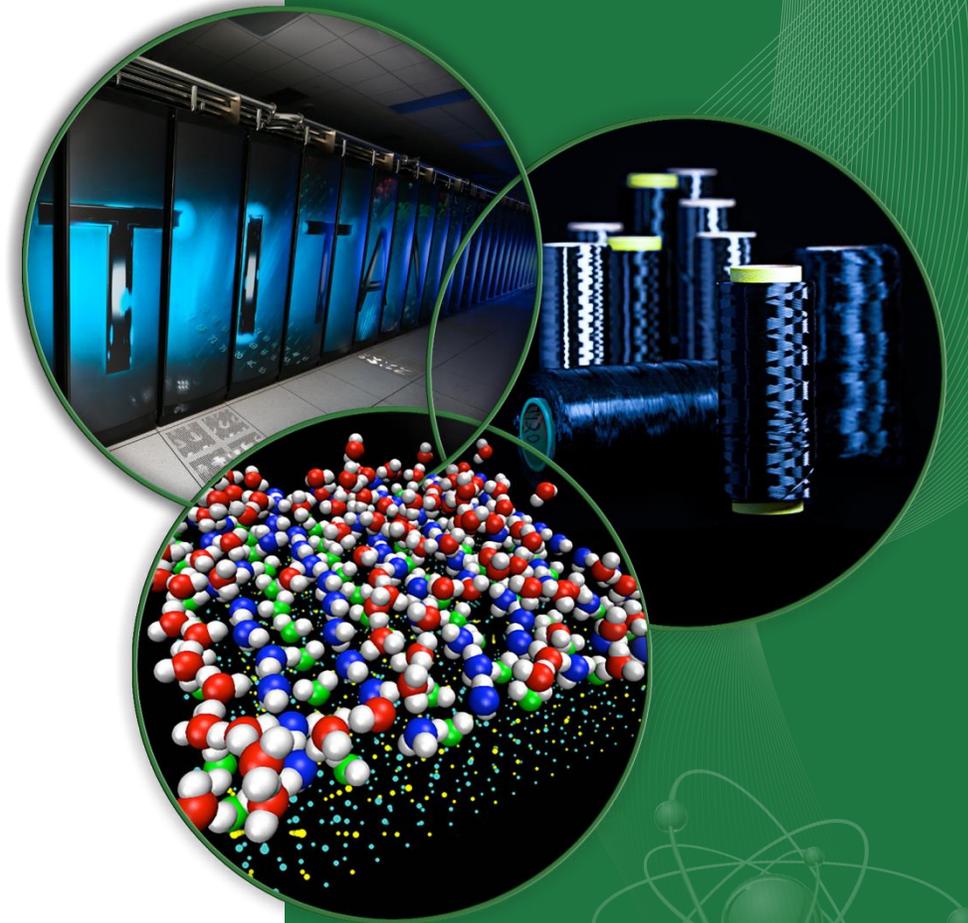
ORNL is managed by UT-Battelle
for the US Department of Energy

Georgia Tech  **College of Computing**
Computational Science and Engineering

<http://ft.ornl.gov> ♦ vetter@computer.org



 **OAK RIDGE**
National Laboratory



Overview

- Our community has major challenges in HPC as we move to extreme scale
 - Power, Performance, Resilience, Productivity
 - New technologies emerging to address some of these challenges
 - Heterogeneous computing
 - Multimode memory systems including nonvolatile memory
 - Not just HPC: Most uncertainty in at least two decades
 - Exascale includes even more diversity and uncertainty
- **We need performance prediction and portability tools now more than ever!**
- Aspen is a tool for structured design and analysis
 - Co-design applications and architectures for performance, power, resiliency
 - Automatic model generation
 - Scalable to distributed scientific workflows
- OpenARC research compiler is a vehicle for
 - Understanding how to automate platform specific optimizations
 - Developing performance portable code

DOE's Office of Science Computation User Facilities



- DOE is leader in open High-Performance Computing
- Provide the world's most powerful computational tools for open science
- Access is free to researchers who publish
- Boost US competitiveness
- Attract the best and brightest researchers



NERSC
Edison is 2.57 PF



ALCF
Mira is 10 PF



OLCF
Titan is 27 PF

HPC Trends



**the 10-Megabyte
Computer System**



**Only
\$5995**
COMPLETE

New From IMSAI®

- 10-Megabyte Hard Disk
- 5¼" Dual-Density Floppy Disk Back-up
- 8-Bit Microprocessor (Optional 16-bit Microprocessor)
- Memory-Mapped Video Display Board
- Disk Controller
- Standard 64K RAM (Optional 256K RAM)
- 10-Slot S-100 Motherboard
- 28-Amp Power Supply
- 12" Monitor
- Standard Intelligent 62-Key ASCII Keyboard (Optional Intelligent 86-Key ASCII Extended Keyboard)
- 132-Column Dot-Matrix Printer
- CP/M* Operating System

**You Read It Right ...
All for \$5995!**

IMSAI® ...Thinking ahead for the 80's

Computer Division of the Fischer-Freitas Corporation
415/635-7615 910 81st Avenue, Bldg. 14 • Oakland, CA 94621

*CP/M is a trademark of Digital Research. Imsai is a trademark of the Fischer-Freitas Corporation

Exascale architecture targets circa 2009

2009 Exascale Challenges Workshop in San Diego

Attendees envisioned two possible architectural swim lanes:

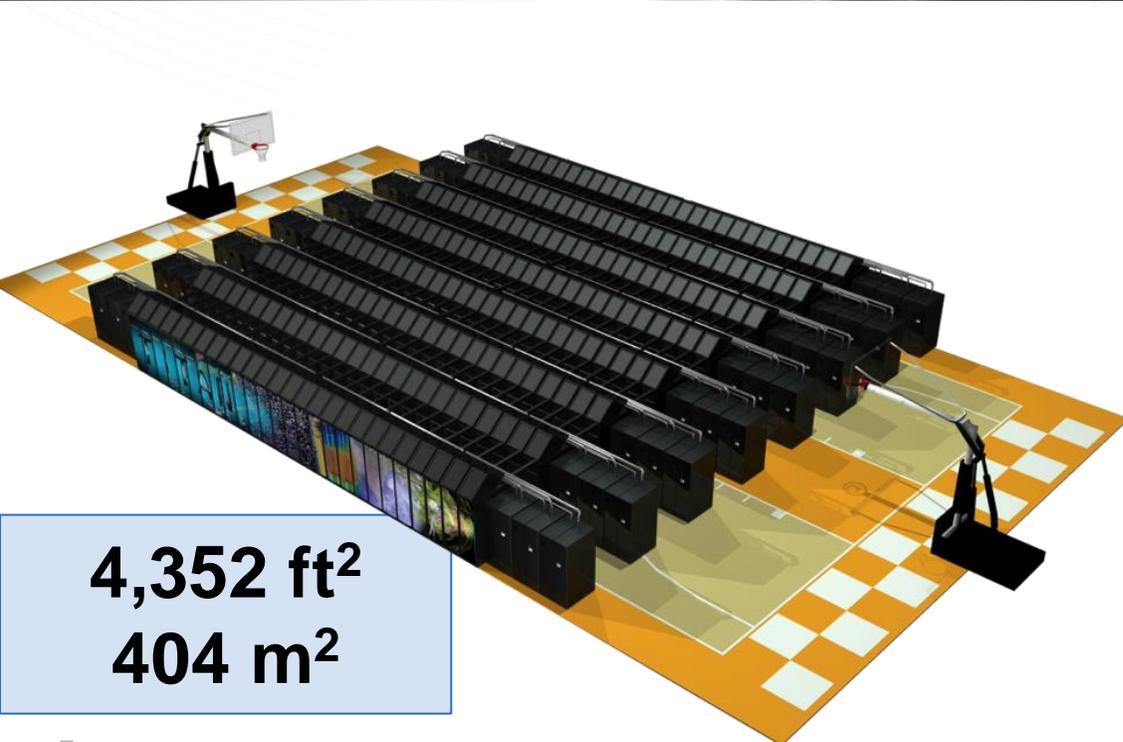
1. Homogeneous many-core thin-node system
2. Heterogeneous (accelerator + CPU) fat-node system

System attributes	2009	“Pre-Exascale”		“Exascale”	
System peak	2 PF	100-200 PF/s		1 Exaflop/s	
Power	6 MW	15 MW		20 MW	
System memory	0.3 PB	5 PB		32–64 PB	
Storage	15 PB	150 PB		500 PB	
Node performance	125 GF	0.5 TF	7 TF	1 TF	10 TF
Node memory BW	25 GB/s	0.1 TB/s	1 TB/s	0.4 TB/s	4 TB/s
Node concurrency	12	O(100)	O(1,000)	O(1,000)	O(10,000)
System size (nodes)	18,700	500,000	50,000	1,000,000	100,000
Node interconnect BW	1.5 GB/s	150 GB/s	1 TB/s	250 GB/s	2 TB/s
IO Bandwidth	0.2 TB/s	10 TB/s		30-60 TB/s	
MTTI	day	O(1 day)		O(0.1 day)	

ASCR Computing At a Glance

System attributes	NERSC Now	OLCF Now	ALCF Now	NERSC Upgrade	OLCF Upgrade	ALCF Upgrades	
Planned Installation	Edison	TITAN	MIRA	Cori 2016	Summit 2017-2018	Theta 2016	Aurora 2018-2019
System peak (PF)	2.6	27	10	> 30	150	>8.5	180
Peak Power (MW)	2	9	4.8	< 3.7	10	1.7	13
Total system memory	357 TB	710TB	768TB	~1 PB DDR4 + High Bandwidth Memory (HBM)+1.5PB persistent memory	> 1.74 PB DDR4 + HBM + 2.8 PB persistent memory	>480 TB DDR4 + High Bandwidth Memory (HBM)	> 7 PB High Bandwidth On-Package Memory Local Memory and Persistent Memory
Node performance (TF)	0.460	1.452	0.204	> 3	> 40	> 3	> 17 times Mira
Node processors	Intel Ivy Bridge	AMD Opteron Nvidia Kepler	64-bit PowerPC A2	Intel Knights Landing many core CPUs Intel Haswell CPU in data partition	Multiple IBM Power9 CPUs & multiple Nvidia Voltas GPUS	Intel Knights Landing Xeon Phi many core CPUs	Knights Hill Xeon Phi many core CPUs
System size (nodes)	5,600 nodes	18,688 nodes	49,152	9,300 nodes 1,900 nodes in data partition	~3,500 nodes	>2,500 nodes	>50,000 nodes
System Interconnect	Aries	Gemini	5D Torus	Aries	Dual Rail EDR-IB	Aries	2 nd Generation Intel Omni-Path Architecture
File System	7.6 PB 168 GB/s, Lustre®	32 PB 1 TB/s, Lustre®	26 PB 300 GB/s GPFS™	28 PB 744 GB/s Lustre®	120 PB 1 TB/s GPFS™	10PB, 210 GB/s Lustre initial	150 PB 1 TB/s Lustre®

ORNL's "Titan" Hybrid System: Cray XK7 with AMD Opteron and NVIDIA Tesla processors



4,352 ft²
404 m²

SYSTEM SPECIFICATIONS:

- Peak performance of 27.1 PF (24.5 & 2.6)
- 18,688 Compute Nodes each with:
- 16-Core AMD Opteron CPU (32 GB)
- NVIDIA Tesla "K20x" GPU (6 GB)
- 512 Service and I/O nodes
- 200 Cabinets
- 710 TB total system memory
- Cray Gemini 3D Torus Interconnect

2017 OLCF Leadership System

Hybrid CPU/GPU architecture

Vendor: **IBM (Prime) / NVIDIA™ / Mellanox Technologies®**

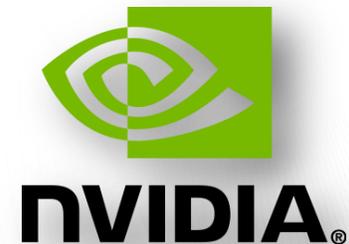
At least 5X Titan's Application Performance

Approximately 3,400 nodes, each with:

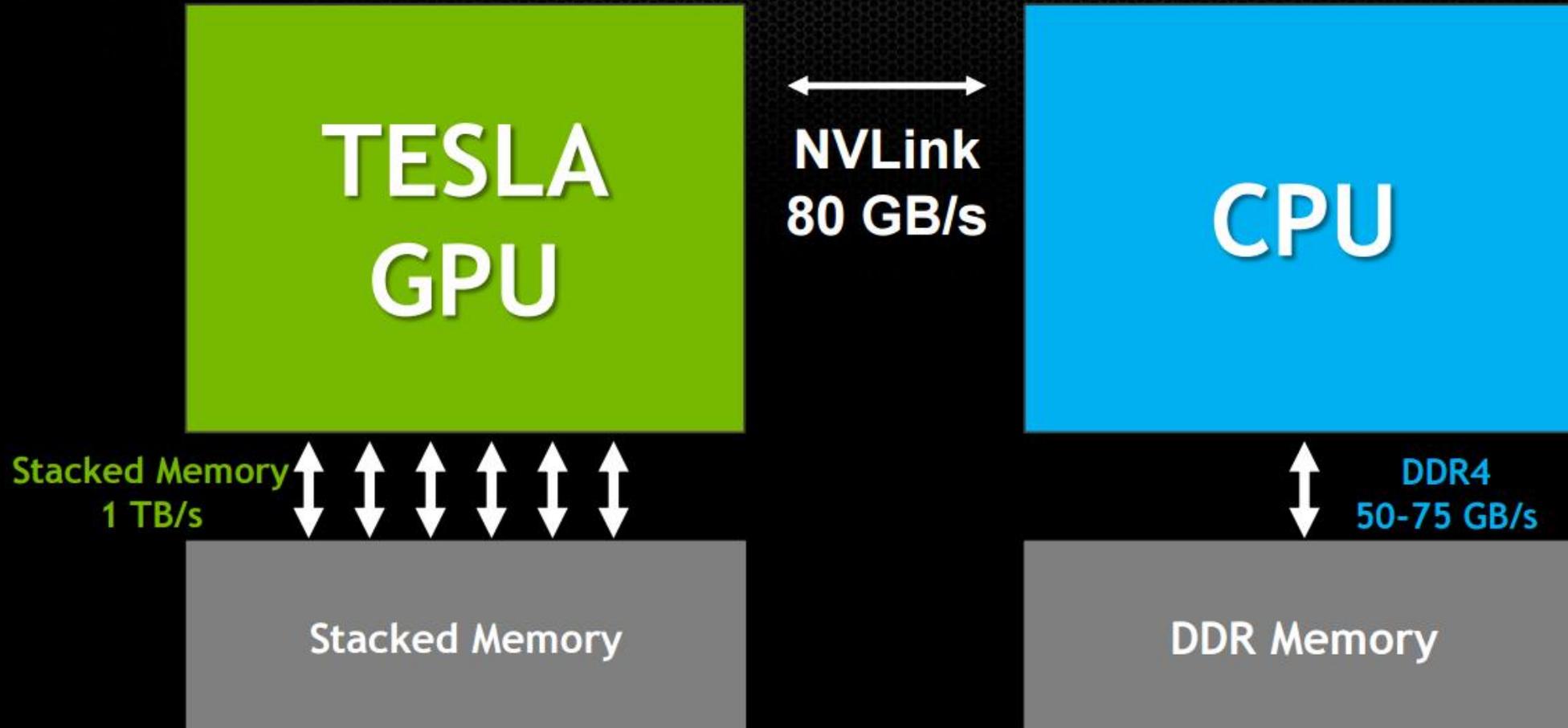
- Multiple IBM POWER9 CPUs and multiple NVIDIA Tesla® GPUs using the NVIDIA Volta architecture
- CPUs and GPUs completely connected with high speed NVLink
- Large coherent memory: over 512 GB (HBM + DDR4)
 - all directly addressable from the CPUs and GPUs
- An additional 800 GB of NVRAM, either a burst buffer or as extended memory
- Over 40 TF peak performance

Dual-rail Mellanox® EDR-IB full, non-blocking fat-tree interconnect

IBM Elastic Storage (GPFS™) - 1TB/s I/O and 120 PB disk capacity.



NVLink Enables Data Transfer At Speed of CPU Memory



OLCF-5 Projections

Date	2009	2012	2017	OLCF-5 2022
System	Jaguar	Titan	Summit	Exascale
System peak	2.3 Peta	27 Peta	150+ Peta	1-2 Exa
System memory	0.3 PB	0.7 PB	2-5 PB	10-20 PB
NVM per node	none	none	800 GB	~2 TB
Storage	15 PB	32 PB	120 PB	~300 PB
MTTI	days	days	days	O(1 day)
Power	7 MW	9 MW	10 MW	~20 MW
Node architecture	CPU 12 core	CPU + GPU	X CPU + Y GPU	X loc + Y toc
System size (nodes)	18,700	18,700	3,400	How fat?
Node performance	125 GF	1.5 TF	40 TF	depends (X,Y)
Node memory BW	25 GB/s	25 - 200 GB/s	100 – 1000 GB/s	10x fast vs slow
Interconnect BW	1.5 GB/s	6.4 GB/s	25 GB/s	4x each gen
IO Bandwidth	0.2 TB/s	1 TB/s	1 TB/s	flat

International Progress: Japan

Feature and configuration of FX100



Tofu interconnect 2

- 12.5 GB/s×2 (in/out)/link
- 10 links/node
- Optical technology

CPU Memory Board

- Three CPUs
- 3 x 8 Micron's HMCs
- 8 opt modules, for inter-chassis connections

Cabinet

- Up to 216 nodes/cabinet
- High-density
- 100% water cooled with EXCU (option)



Fujitsu designed SPARC64 XIfx

- 1TF~ (DP)/2TF~ (SP)
- 32 + 2 core CPU
- HPC-ACE2 support
- Tofu2 integrated

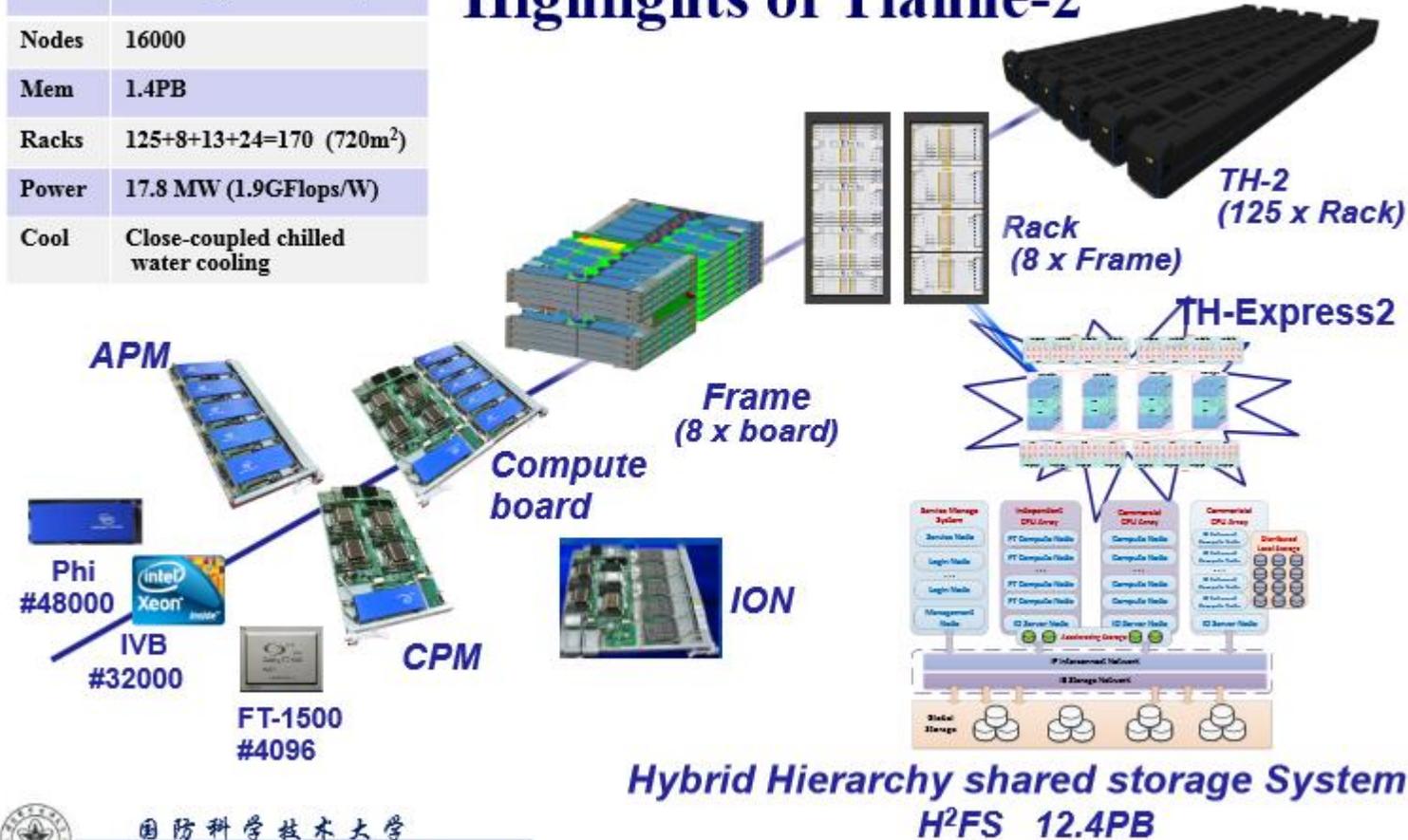
Chassis

- 1 CPU/1 node
- 12 nodes/2U Chassis
- Water cooled

International Progress: China

Perf	54.9PFlops / 33.86PFlops
Nodes	16000
Mem	1.4PB
Racks	125+8+13+24=170 (720m ²)
Power	17.8 MW (1.9GFlops/W)
Cool	Close-coupled chilled water cooling

Highlights of Tianhe-2



China May Develop Two 100 Petaflop Machines Within a Year

August 26, 2015 by [Rich Brueckner](#) 1 Comment

In this special guest feature from [Scientific Computing World](#), Tom Wilkie looks at the next generation of supercomputers coming to China.



Within the next 12 months, China expects to be operating not one but two 100 Petaflop

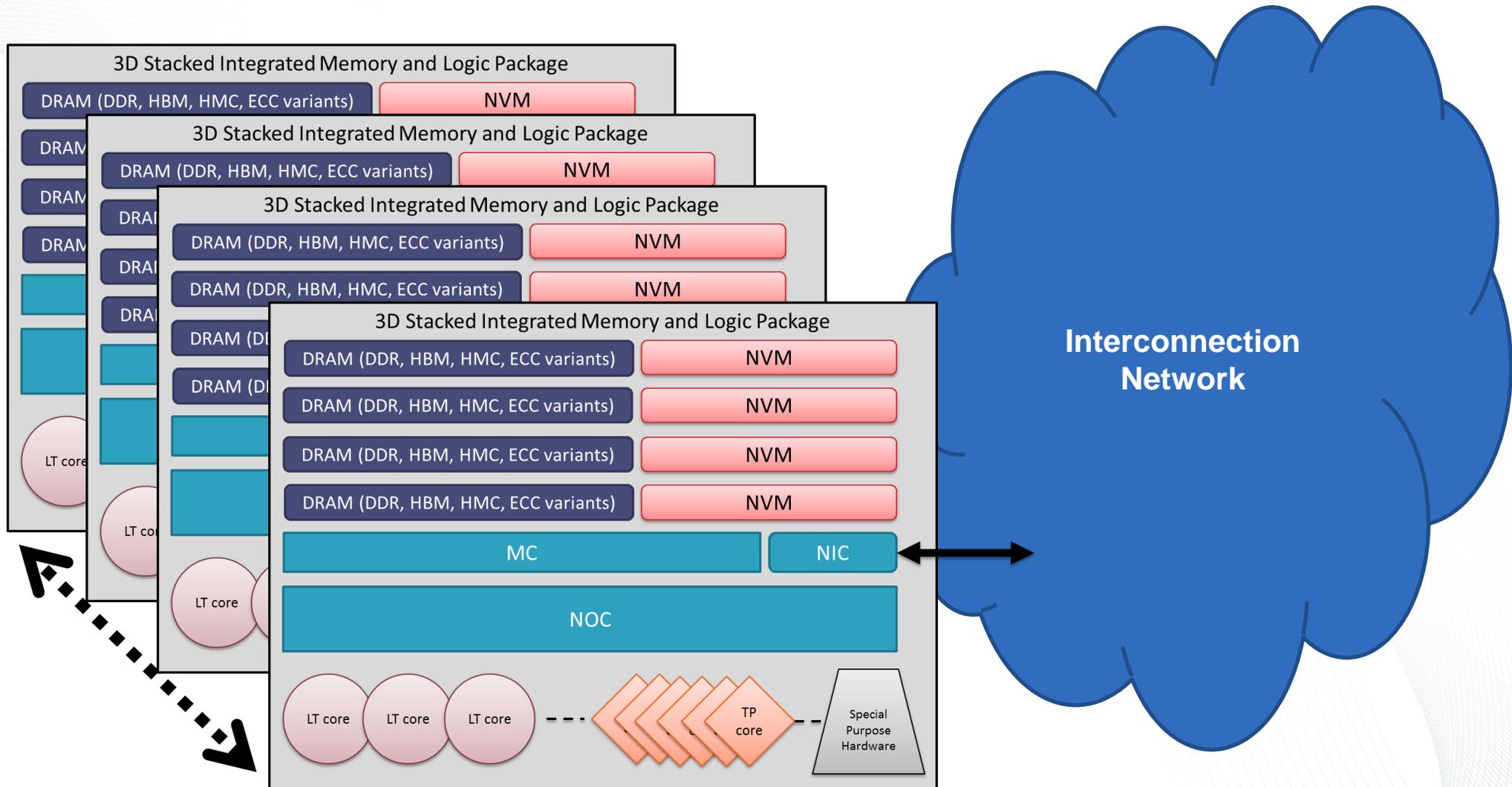
computers, each containing (different) Chinese-made processors, and both coming online about a year before the United States' 100 Petaflop machines being developed under the [Coral initiative](#).

Ironically, the CPU for one machine appears very similar to a technology abandoned by the USA in 2007, and the US Government, through its export embargo, has encouraged China to develop its own accelerator for the other machine.



国防科学技术大学
National University of Defense Technology

Notional Future Architecture

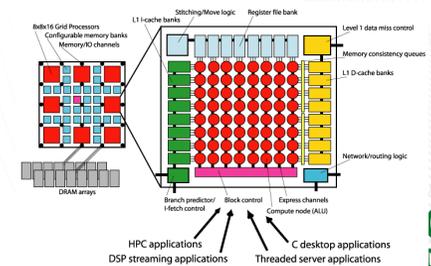
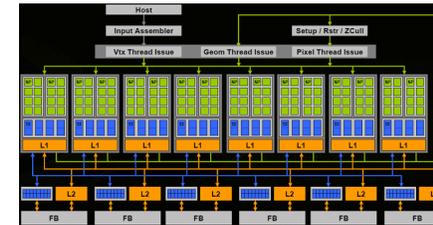
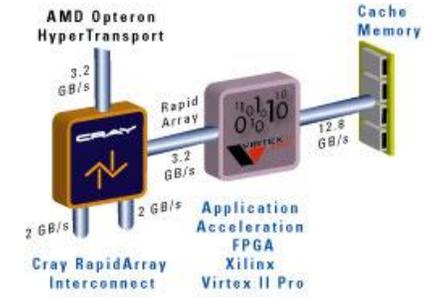
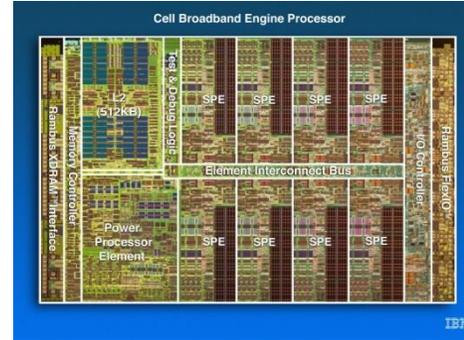


Heterogeneous computing is
here to stay

Earlier Experimental Computing Systems (past decade)

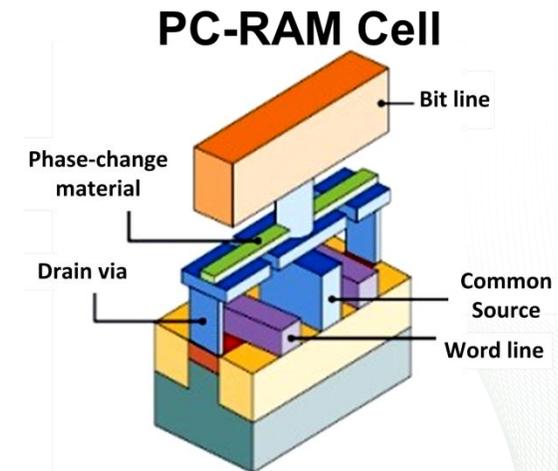
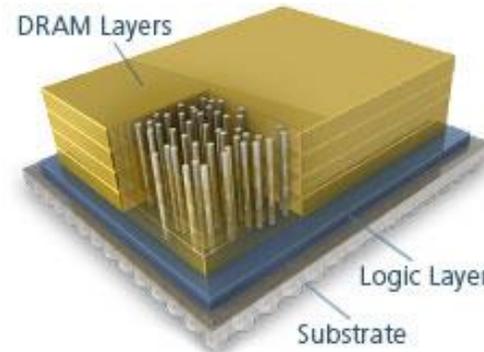
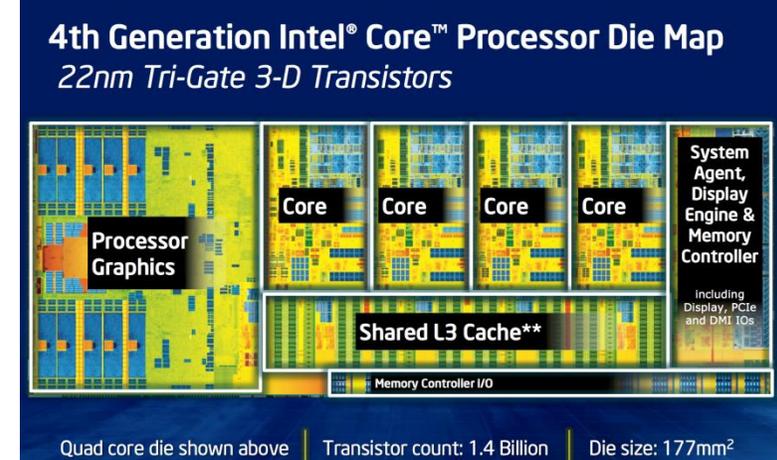
- The past decade has started the trend away from traditional 'simple' architectures
- Examples
 - Cell, GPUs, FPGAs, SoCs, etc
- Lessons learned?
- Mainly driven by facilities costs and successful (sometimes heroic) application examples

Popular architectures since ~2004



Emerging Computing Architectures – Future

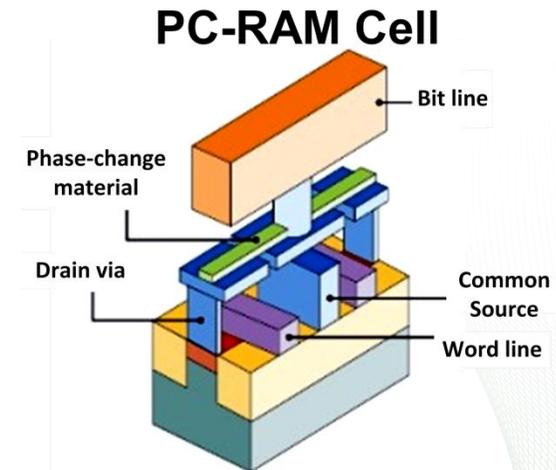
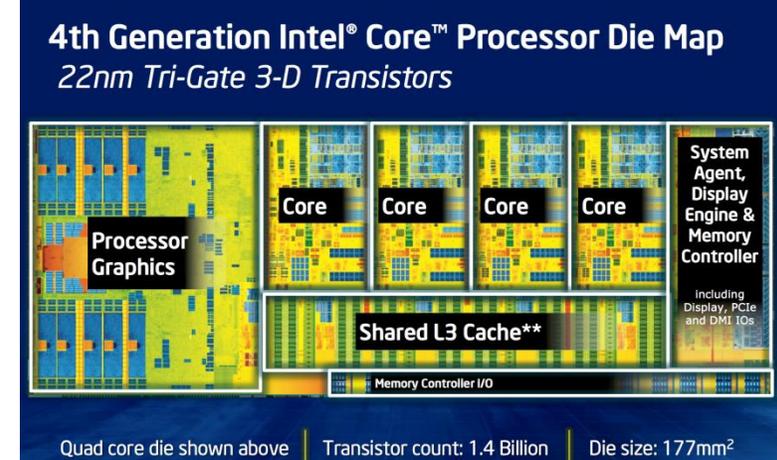
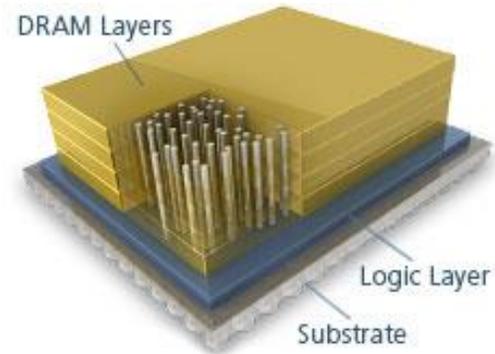
- Heterogeneous processing
 - Latency tolerant cores
 - Throughput cores
 - Special purpose hardware (e.g., AES, MPEG, RND)
 - Fused, configurable memory
- Memory
 - 2.5D and 3D Stacking
 - HMC, HBM, WIDEIO₂, LPDDR₄, etc
 - New devices (PCRAM, ReRAM)
- Interconnects
 - Collective offload
 - Scalable topologies
- Storage
 - Active storage
 - Non-traditional storage architectures (key-value stores)
- Improving performance and programmability in face of increasing complexity
 - Power, resilience



HPC (mobile, enterprise, embedded) computer design is more fluid now than in the past two decades.

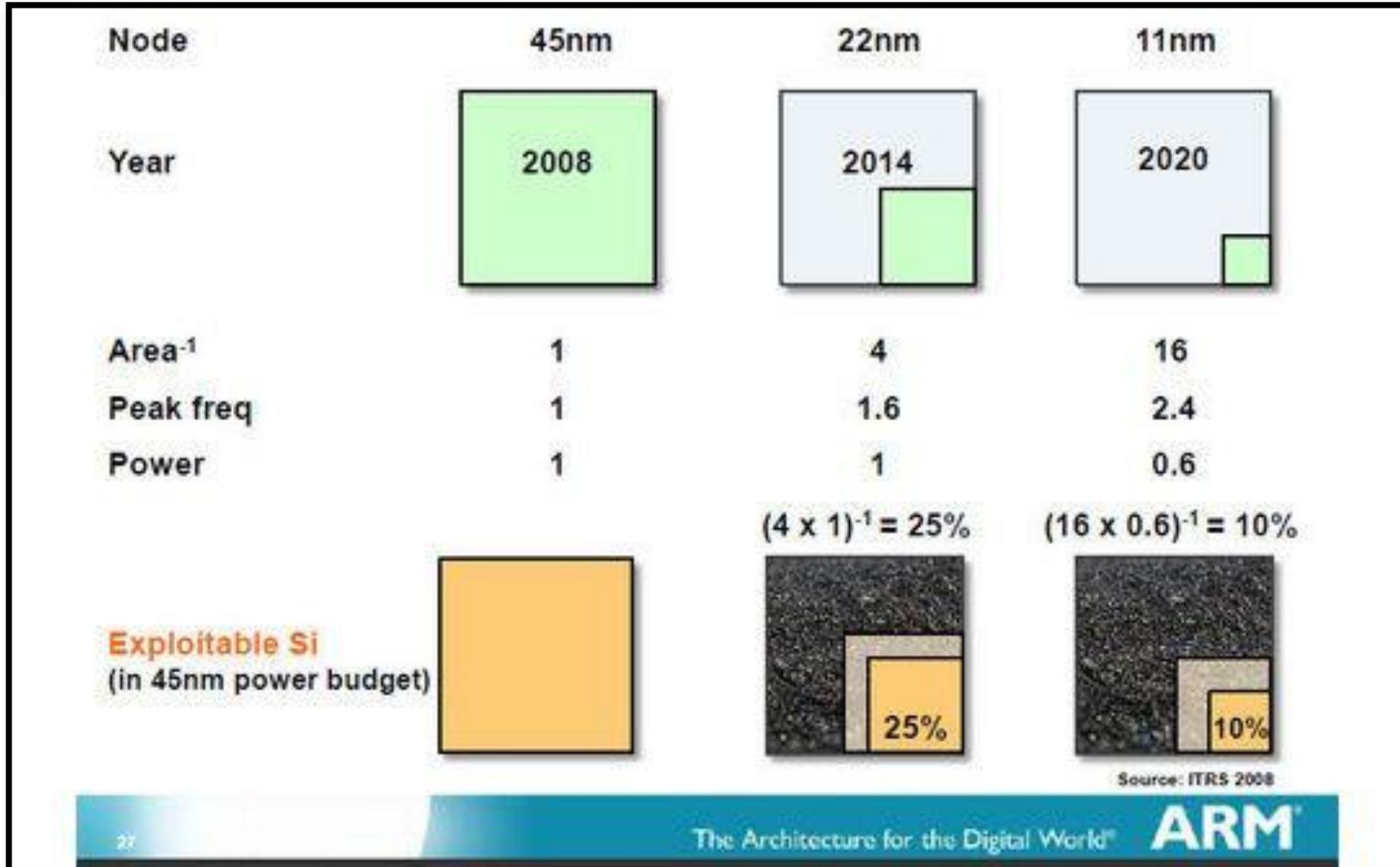
Emerging Computing Architectures – Future

- Heterogeneous processing
 - Latency tolerant cores
 - Throughput cores
 - Special purpose hardware (e.g., AES, MPEG, RND)
 - Fused, configurable memory
- Memory
 - 2.5D and 3D Stacking
 - HMC, HBM, WIDEIO₂, LPDDR₄, etc
 - New devices (PCRAM, ReRAM)
- Interconnects
 - Collective offload
 - Scalable topologies
- Storage
 - Active storage
 - Non-traditional storage architectures (key-value stores)
- Improving performance and programmability in face of increasing complexity
 - Power, resilience



HPC (mobile, enterprise, embedded) computer design is more fluid now than in the past two decades.

Dark Silicon Facilitates Heterogeneity and Specialization



Recent announcements

Nvidia and IBM create GPU interconnect for faster supercomputing

"NVLink" shares up to 80GB of data per second between CPUs and GPUs.

It Begins: AMD Announces Its First ARM Based Server SoC, 64-bit/8-core Opteron A1100

by Anand Lal Shimpi on January 28, 2014 6:35 PM EST

Posted in CPUs IT Computing Enterprise enterprise CPUs AMD Opteron Opteron A1100 ARM

Intel's 14nm Broadwell GPU takes shape, indicates major improvements over Haswell

By Sebastian Anthony on November 5, 2013 at 10:21 am | 16 Comments



"SEATTLE" 64-BIT ARM SERVER PROCESSOR FIRST 28NM ARM Nvidia Jetson TK1 mini supercomputer is up for pre-order

Will ship on 15 May
By Lee Bell
Fri May 02 2014, 11:00 AM

- Indus prove
- The
- Serv
- CPU
- 2-4x impr
- 8 cor
- Base
- Ext
- CPU
- Serv
- Leg
- Stor

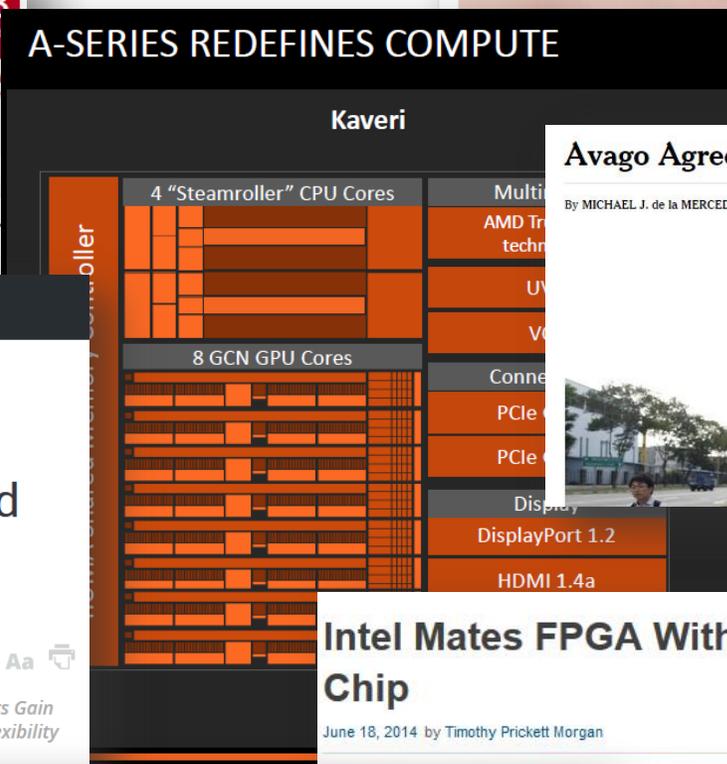
MarketWatch

Altera and IBM Unveil FPGA-accelerated POWER Systems with Coherent Shared Memory

By
Published: Nov 17, 2014 8:00 a.m. ET

POWER8 Systems that Leverage Reprogrammable FPGA Accelerators Gain Significant Improvements in System Performance, Efficiency and Flexibility

NEW ORLEANS, Nov. 17, 2014 /PRNewswire/ Corporation ALTR, +0.00% and IBM IBM, +0.00% FPGA-based acceleration platform that cohes CPU leveraging IBM's Coherent Accelerator I reconfigurable hardware accelerator feature FPGA and processor which significantly imp flexibility in high-performance computing (H and IBM are presenting several POWER8 sys using FPGAs at SuperComputing 2014.



Avago Agrees to Buy Broadcom for \$37 Billion

By MICHAEL J. de la MERCED and CHAD BRAY MAY 28, 2015



started open-sourcing the Linux driver for Broadwell's GPU. Broadwell is the 14nm die shrink of Intel's

Intel Mates FPGA With Future Xeon Server Chip

June 18, 2014 by Timothy Prickett Morgan

Intel is taking field programmable gate arrays seriously as a means of accelerating applications and has crafted a hybrid chip that marries an FPGA to a Xeon E5 processor and puts them in the same processor socket.

Intel to acquire Altera for \$54 a share

Monday, 1 Jun 2015 | 8:33 AM ET



Working together through the OpenPOWER

Nvidia CEO Jen-Hsun Huang described the Jetson TK1 as a modified GPU, and the

While the CPU expected to Broadwell's GPU looks like the Iris GPU a modified GPU, and the

Intel is taking field programmable gate arrays seriously as a means of accelerating applications and has crafted a hybrid chip that marries an FPGA to a Xeon E5 processor and puts them in the same processor socket.

Tighter Integration will expand workload possibilities

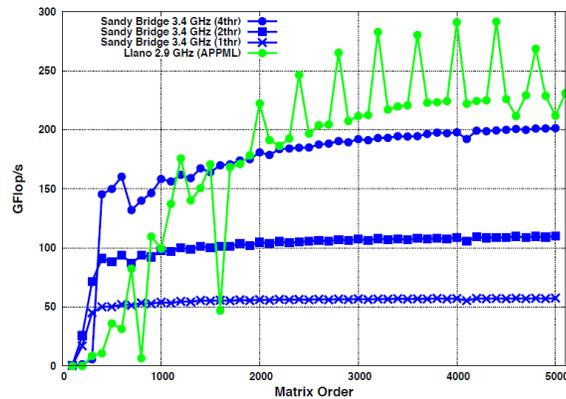
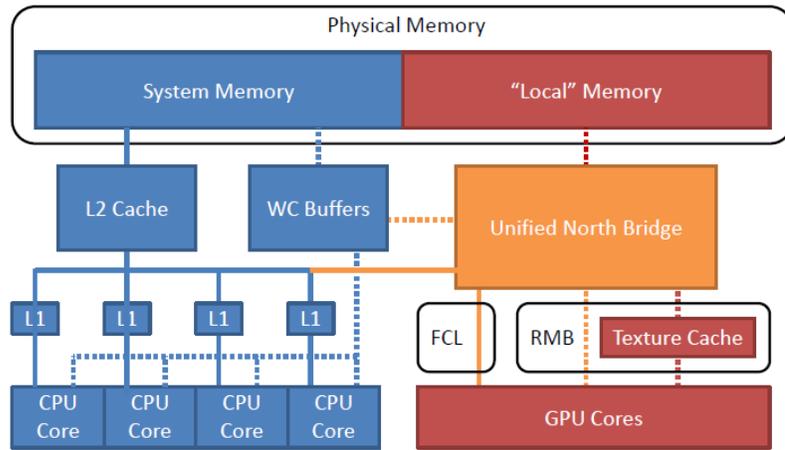
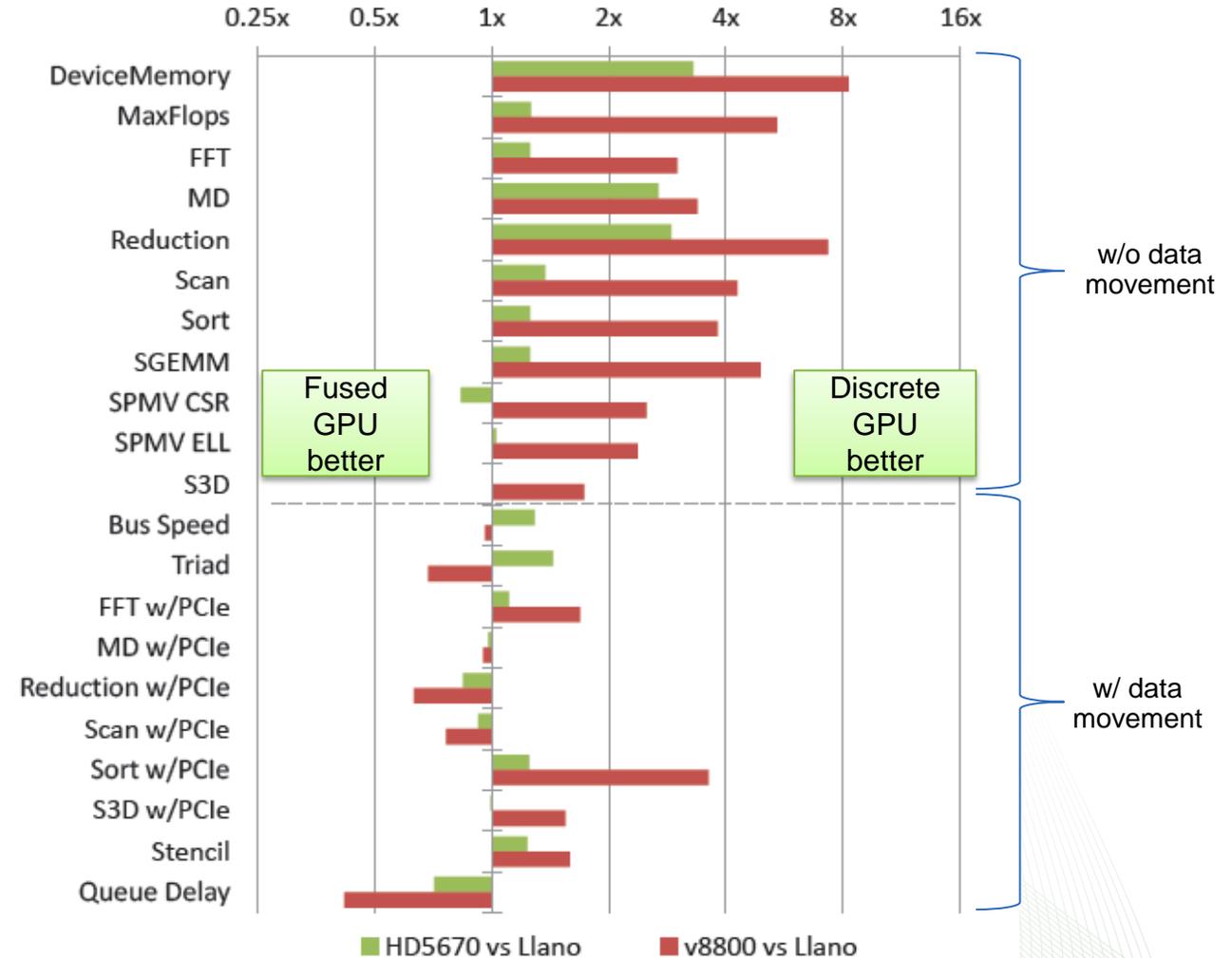
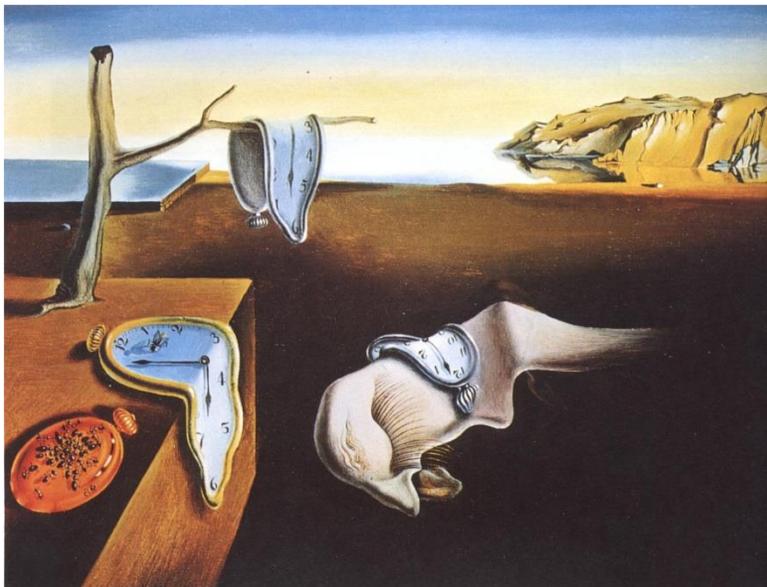


Figure 3: SGEMM Performance (one, two, and four CPU threads for Sandy Bridge and the OpenCL-based AMD APPML for Llano's fGPU)



New and Improved Memory Systems are the Next Big Thing



Emerging Computing Architectures – Future

- Heterogeneous processing
 - Latency tolerant cores
 - Throughput cores
 - Special purpose hardware (e.g., AES, MPEG, RND)
 - Fused, configurable memory

Memory

- 2.5D and 3D Stacking
- HMC, HBM, WIDEIO₂, LPDDR₄, etc
- New devices (PCRAM, ReRAM)
- New interfaces

Interconnects

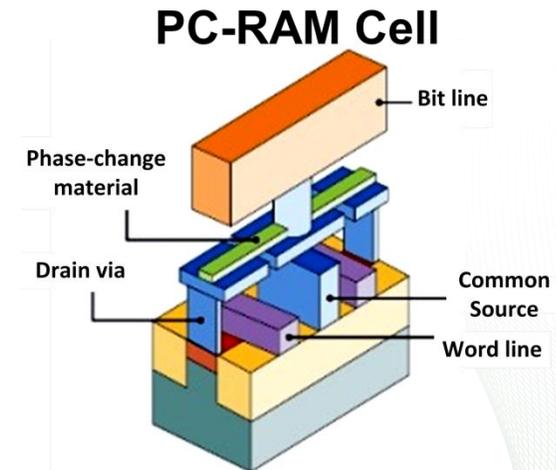
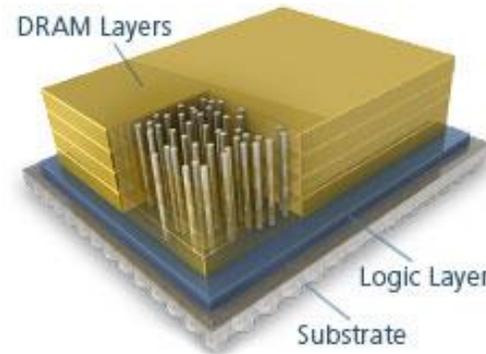
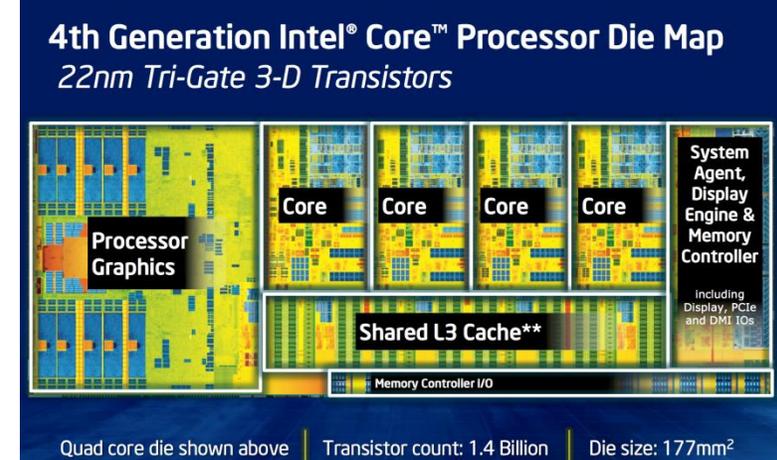
- Collective offload
- Scalable topologies

Storage

- Active storage
- Non-traditional storage architectures (key-value stores)

Improving performance and programmability in face of increasing complexity

- Power, resilience



HPC (mobile, enterprise, embedded) computer design is more fluid now than in the past two decades.

Exascale architecture targets

defined at 2009 Exascale Challenges Workshop in San Diego

Where we are going “off the tracks” is data movement between nodes and from node to storage

Summit: Interconnect BW= 25 GB/s, I/O BW= 1 TB/s

System attributes	2009	“Pre-Exascale”		“Exascale”	
System peak	2 PF	100-200 PF/s		1 Exaflop/s	
Power	6 MW	15 MW		20 MW	
System memory	0.3 PB	5 PB		32–64 PB	
Storage	15 PB	150 PB		500 PB	
Node performance	125 GF	0.5 TF	7 TF	1 TF	10 TF
Node memory BW	25 GB/s	0.1 TB/s	1 TB/s	0.4 TB/s	4 TB/s
Node concurrency	12	O(100)	O(1,000)	O(1,000)	O(10,000)
System size (nodes)	18,700	500,000	50,000	1,000,000	100,000
Node interconnect BW	1.5 GB/s	150 GB/s	1 TB/s	250 GB/s	2 TB/s
IO Bandwidth	0.2 TB/s	10 TB/s		30-60 TB/s	
MTTI	day	O(1 day)		O(0.1 day)	

Exascale architecture targets

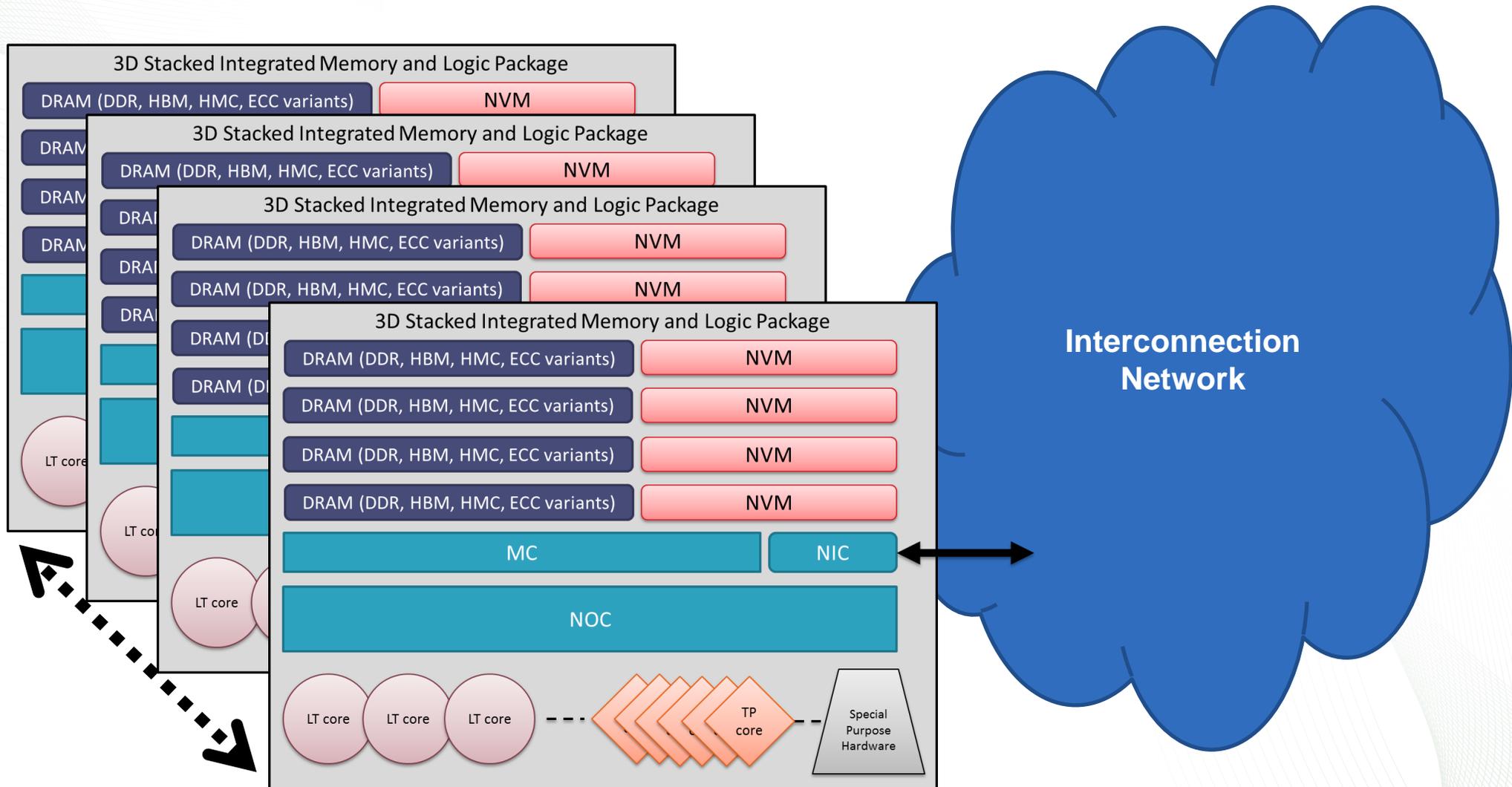
defined at 2009 Exascale Challenges Workshop in San Diego

Where we are going “off the tracks” is data capacity, and movement between nodes and from node to storage

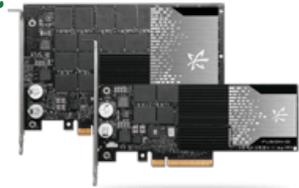
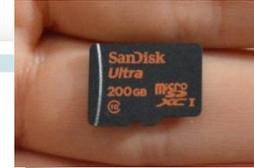
Summit: Interconnect BW= 25 GB/s, I/O BW= 1 TB/s

System attributes	2009	“Pre-Exascale”		“Exascale”	
System peak	2 PF	100-200 PF/s		1 Exaflop/s	
Power	6 MW	15 MW		20 MW	
System memory	0.3 PB	5 PB		32–64 PB	
Storage	15 PB	150 PB		500 PB	
Node performance	125 GF	0.5 TF	7 TF	1 TF	10 TF
Node memory BW	25 GB/s	0.1 TB/s	1 TB/s	0.4 TB/s	4 TB/s
Node concurrency	12	O(100)	O(1,000)	O(1,000)	O(10,000)
System size (nodes)	18,700	500,000	50,000	1,000,000	100,000
Node interconnect BW	1.5 GB/s	150 GB/s	1 TB/s	250 GB/s	2 TB/s
IO Bandwidth	0.2 TB/s	10 TB/s		30-60 TB/s	
MTTI	day	O(1 day)		O(0.1 day)	

Notional Future Architecture



NVRAM Technology Continues to Improve – Driven by Market Forces



designlines MEMORY

News & Analysis

3D NAND Production Starts at Samsung

Peter Clarke

8/6/2013 08:05 AM EDT
16 comments

Like 17 Tweet 7

LONDON — Samsung Electronics has started production of a 128 Gbit NAND memory chip with multiple layers, and claims it is the world's first.

The memory is based on conventional floating-gate technology. In the vertical arrangement, the reliability between a factored conventional floating-gate and a V-NAND is a press release.

The technology is capable of stacking up to 24 layers, but Samsung did not disclose how many layers it had used in its 128 Gbit vertical NAND, nor whether the memory cells are multilevel cell or whether it had relaxed the design geometry from the leading edge in 2D memory, which stands at about 19 or 16 nm.

The company did say that the memory would provide improvements in performance and area ratio, and a V-NAND chip is suitable for a wide range of consumer and commercial applications including embedded NAND storage and solid-state drives.

The V-NAND component has the same memory capacity as a 128

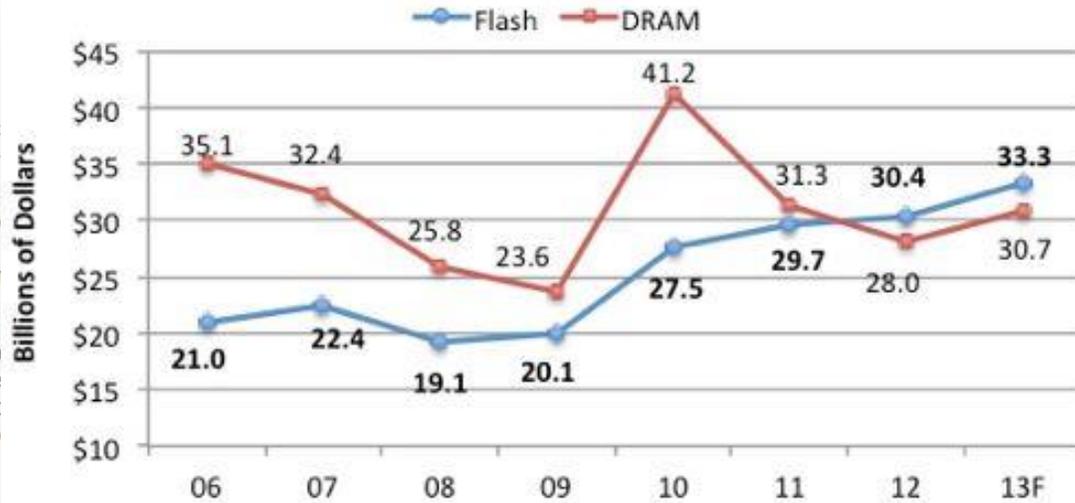
designlines MEMORY

News & Analysis

3D NAND Transition: 15nm Process Technology Takes Shape

Gary Hilson

NO RATINGS
LOGIN TO RATE



memory business unit. Toshiba's 15nm process works in conjunction with improved peripheral circuitry technology chips that achieve the same write speed as chips formed with second generation 19nm process technology, but boost transfer rate to 533 megabits a second -- 1.3 times faster -- employing a high-speed interface.

Nelson said there is room to advance floating gates before moving

http://www.eetasia.com/STATIC/ARTICLE_IMAGES/201212/EEOL_2012DEC28_STOR_MFG_NT_01.jpg

Original URL: http://www.theregister.co.uk/2013/11/01/hp_memristor_2018/

HP 100TB Memristor drives by 2018 – if you're lucky, admits tech titan

Universal memory slow in coming

By Chris Mellor

8+1

Forbes / Tech

JUL 28, 2015 @ 2:46 PM 7,391 VIEWS

Intel And Micron Jointly Announce Game-Changing 3D XPoint Memory Technology

China's Tsinghua Unigroup plans \$23B bid for Micron Technology

CNBC.com staff | @CNBC

Monday, 13 Jul 2015 | 8:41 PM ET



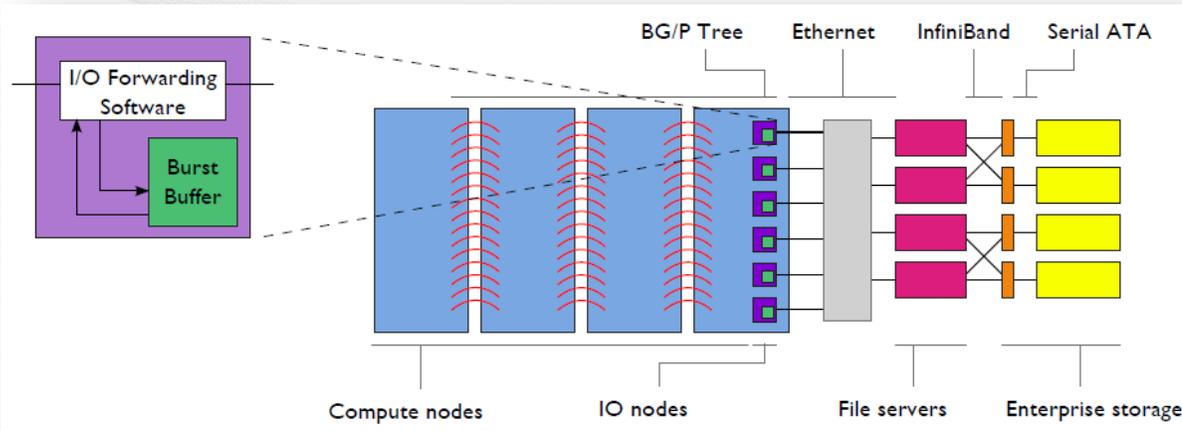
Comparison of emerging memory technologies

	SRAM	DRAM	eDRAM	2D NAND Flash	3D NAND Flash	PCRAM	STTRAM	2D ReRAM	3D ReRAM
Data Retention	N	N	N	Y	Y	Y	Y	Y	Y
Cell Size (F ²)	50-200	4-6	19-26	2-5	<1	4-10	8-40	4	<1
Minimum F demonstrated (nm)	14	25	22	16	64	20	28	27	24
Read Time (ns)	< 1	30	5	10 ⁴	10 ⁴	10-50	3-10	10-50	10-50
Write Time (ns)	< 1	50	5	10 ⁵	10 ⁵	100-300	3-10	10-50	10-50
Number of Rewrites	10 ¹⁶	10 ¹⁶	10 ¹⁶	10 ⁴ -10 ⁵	10 ⁴ -10 ⁵	10 ⁸ -10 ¹⁰	10 ¹⁵	10 ⁸ -10 ¹²	10 ⁸ -10 ¹²
Read Power	Low	Low	Low	High	High	Low	Medium	Medium	Medium
Write Power	Low	Low	Low	High	High	High	Medium	Medium	Medium
Power (other than R/W)	Leakage	Refresh	Refresh	None	None	None	None	Sneak	Sneak
Maturity									

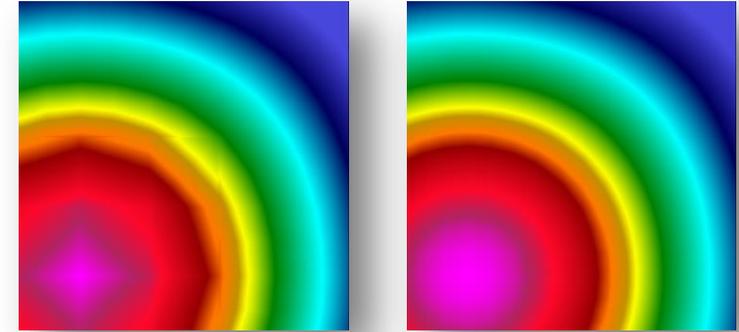
Opportunities for NVM in Emerging Systems

- Burst Buffers

[Liu, et al., MSST 2012]



- In situ visualization



<http://ft.ornl.gov/eavl>

- In-mem tables

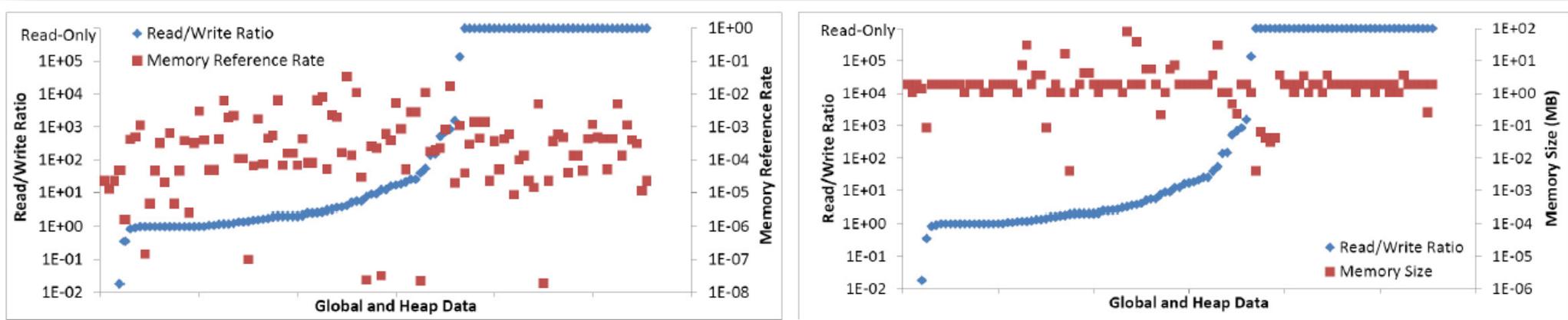


Figure 3: Read/write ratios, memory reference rates and memory object sizes for memory objects in Nek5000

Architectural Uncertainty Demands Methods for Performance Prediction and Portability

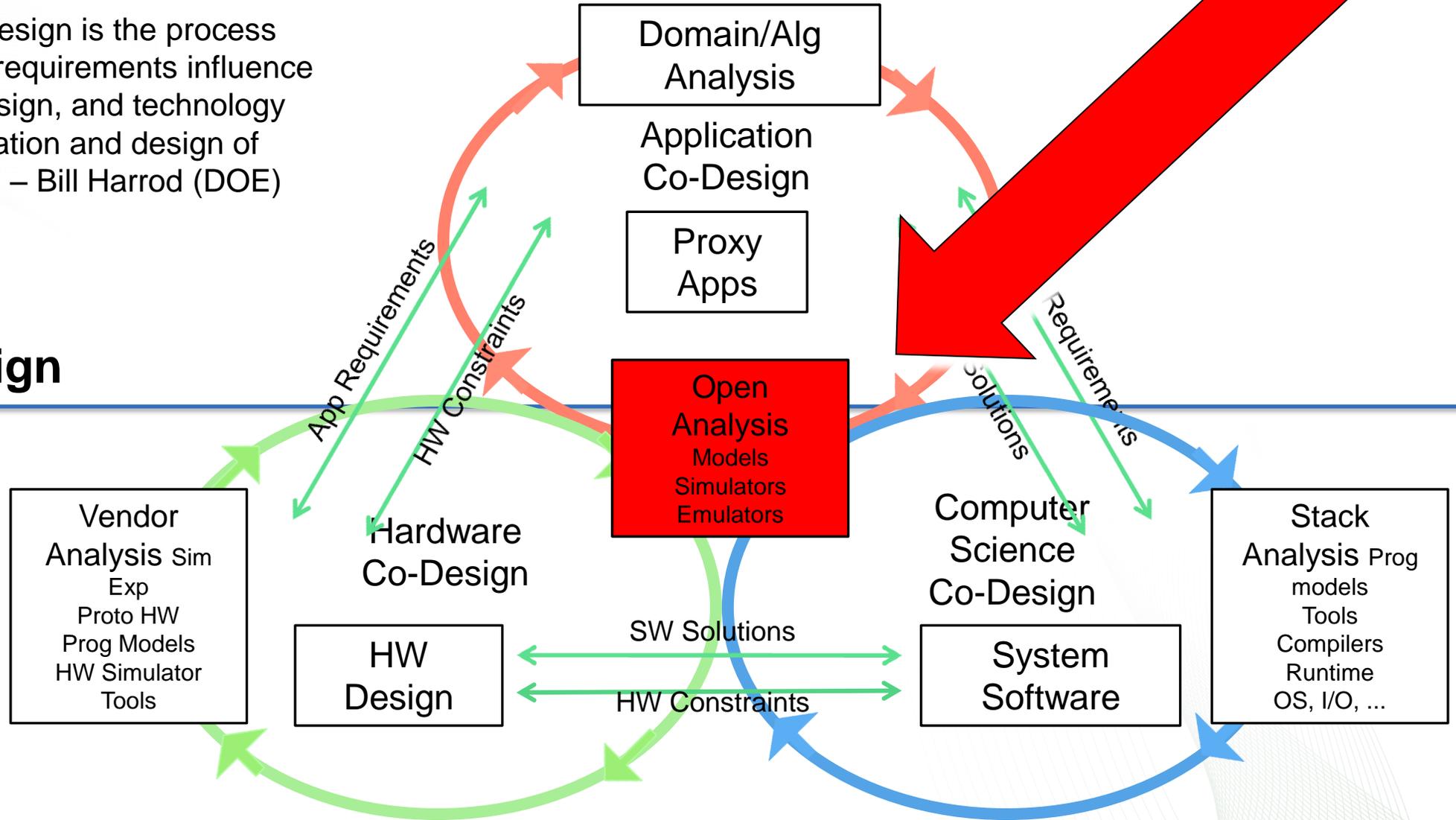
- Performance Prediction – Aspen
- Performance Portability – OpenARC

Workflow within the Exascale Ecosystem

“(Application driven) co-design is the process where scientific problem requirements influence computer architecture design, and technology constraints inform formulation and design of algorithms and software.” – Bill Harrod (DOE)

Application Design

System Design



Prediction Techniques Ranked

	Speed	Ease	Flexibility	Accuracy	Scalability
Ad-hoc Analytical Models	1	3	2	4	1
Structured Analytical Models	1	2	1	4	1
Simulation – Functional	3	2	2	3	3
Simulation – Cycle Accurate	4	2	2	2	4
Hardware Emulation (FPGA)	3	3	3	2	3
Similar hardware measurement	2	1	4	2	2
Node Prototype	2	1	4	1	4
Prototype at Scale	2	1	4	1	2
Final System	-	-	-	-	-

Prediction Techniques Ranked

	Speed	Ease	Flexibility	Accuracy	Scalability
Ad-hoc Analytical Models	1	3	2	4	1
Structured Analytical Models	1	2	1	4	1
<i>Aspen</i>	1	1	1	4	1
Simulation – Functional	3	2	2	3	3
Simulation – Cycle Accurate	4	2	2	2	4
Hardware Emulation (FPGA)	3	3	3	2	3
Similar hardware measurement	2	1	4	2	2
Node Prototype	2	1	4	1	4
Prototype at Scale	2	1	4	1	2
Final System	-	-	-	-	-

Aspen: Abstract Scalable Performance Engineering Notation

Model Creation

- Static analysis via compiler, tools
- Empirical, Historical
- Manual (for future applications)

- ## Representation in Aspen
- Modular
 - Sharable
 - Composable
 - Reflects prog structure



Model Uses

- Interactive tools for graphs, queries
- Design space exploration
- Workload Generation
- Feedback to Runtime Systems

E.g., MD, UHPC CP 1, Lulesh, 3D FFT, CoMD, VPFIT, ...

Source code

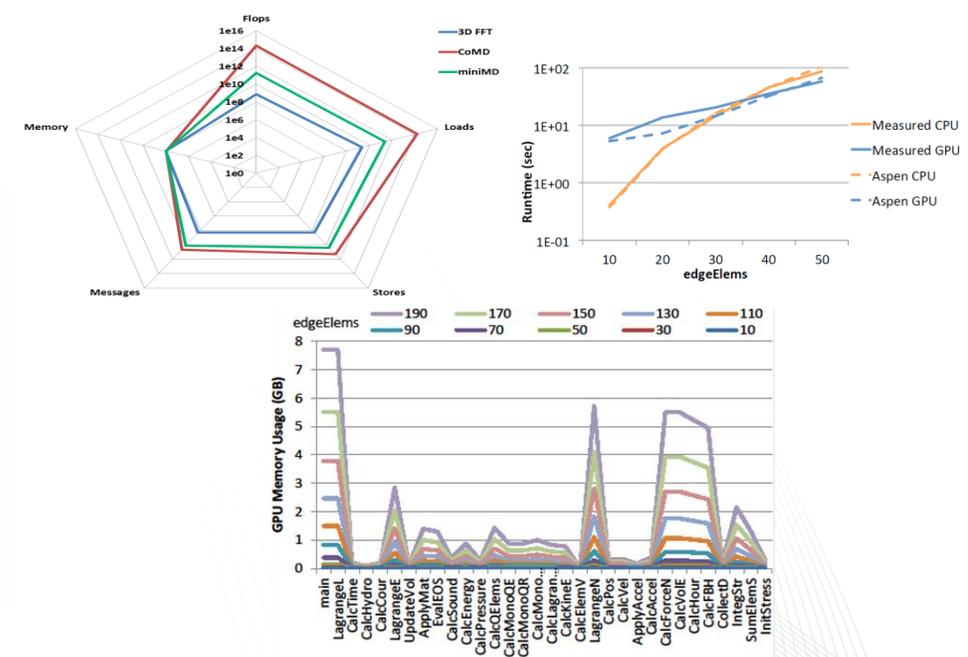
```

2324 static inline
2325 void CalcMonotonicGradientsForElems(Index_t p_nodelist[T_NUMELEM],
2326 Real_t p_x[T_NUMNODE], Real_t p_y[T_NUMNODE], Real_t p_z[T_NUMNODE],
2327 Real_t p_xd[T_NUMNODE], Real_t p_yd[T_NUMNODE], Real_t p_zd[T_NUMNODE],
2328 Real_t p_volo[T_NUMELEM], Real_t p_vnew[T_NUMELEM],
2329 Real_t p_delx_eta[T_NUMELEM], Real_t p_delv_eta[T_NUMELEM],
2330 Real_t p_delx_xi[T_NUMELEM], Real_t p_delv_xi[T_NUMELEM],
2331 Real_t p_delx_eta[T_NUMELEM], Real_t p_delv_eta[T_NUMELEM])
2332 {
2333     Index_t i;
2334     Index_t numElem = m_numElem;
2335     #pragma acc parallel loop independent present(p_vnew, p_nodelist, p_x, p_y, p_z, p_xd, \
2336 p_yd, p_zd, p_volo, p_delx_xi, p_delx_eta, p_delx_eta, p_delv_xi, p_delv_eta, \
2337 p_delv_eta)
2338     for (i = 0; i < numElem; ++i) {
2339         const Real_t ptiny = 1.e-36;
2340         Real_t ax, ay, az;
2341         Real_t dxv, dyv, dzv;
2342
2343         const Index_t *elemToNode = &p_nodelist[*i];
2344         Index_t n0 = elemToNode[0];
2345         Index_t n1 = elemToNode[1];
2346         Index_t n2 = elemToNode[2];
2347         Index_t n3 = elemToNode[3];
2348         Index_t n4 = elemToNode[4];
2349         Index_t n5 = elemToNode[5];
2350         Index_t n6 = elemToNode[6];
2351         Index_t n7 = elemToNode[7];
2352
2353         Real_t x0 = p_x[n0];
    
```

Aspen code

```

147 kernel CalcMonotonicGradients {
148     execute [numElems]
149     {
150         loads [8 * indexWordSize] from nodelist
151         // Load and cache position and velocity.
152         loads/caching [8 * wordSize] from x
153         loads/caching [8 * wordSize] from y
154         loads/caching [8 * wordSize] from z
155
156         loads/caching [8 * wordSize] from xvel
157         loads/caching [8 * wordSize] from yvel
158         loads/caching [8 * wordSize] from zvel
159
160         loads [wordSize] from volo
161         loads [wordSize] from vnew
162         // dx, dy, etc.
163         flops [90] as dp, simd
164         // delvk delix
165         flops [9 * 8 + 3 + 30 + 5] as dp, simd
166         stores [wordSize] to delv_eta
167         // delvi delvi
168         flops [9 * 8 + 3 + 30 + 5] as dp, simd
169         stores [wordSize] to delx_xi
170         // delvj and delvj
171         flops [9 * 8 + 3 + 30 + 5] as dp, simd
172         stores [wordSize] to delv_eta
173     }
174 }
    
```



Researchers are using Aspen for parallel applications, scientific workflows, capacity planning, power, quantum computing, etc



Manual Example of LULESH

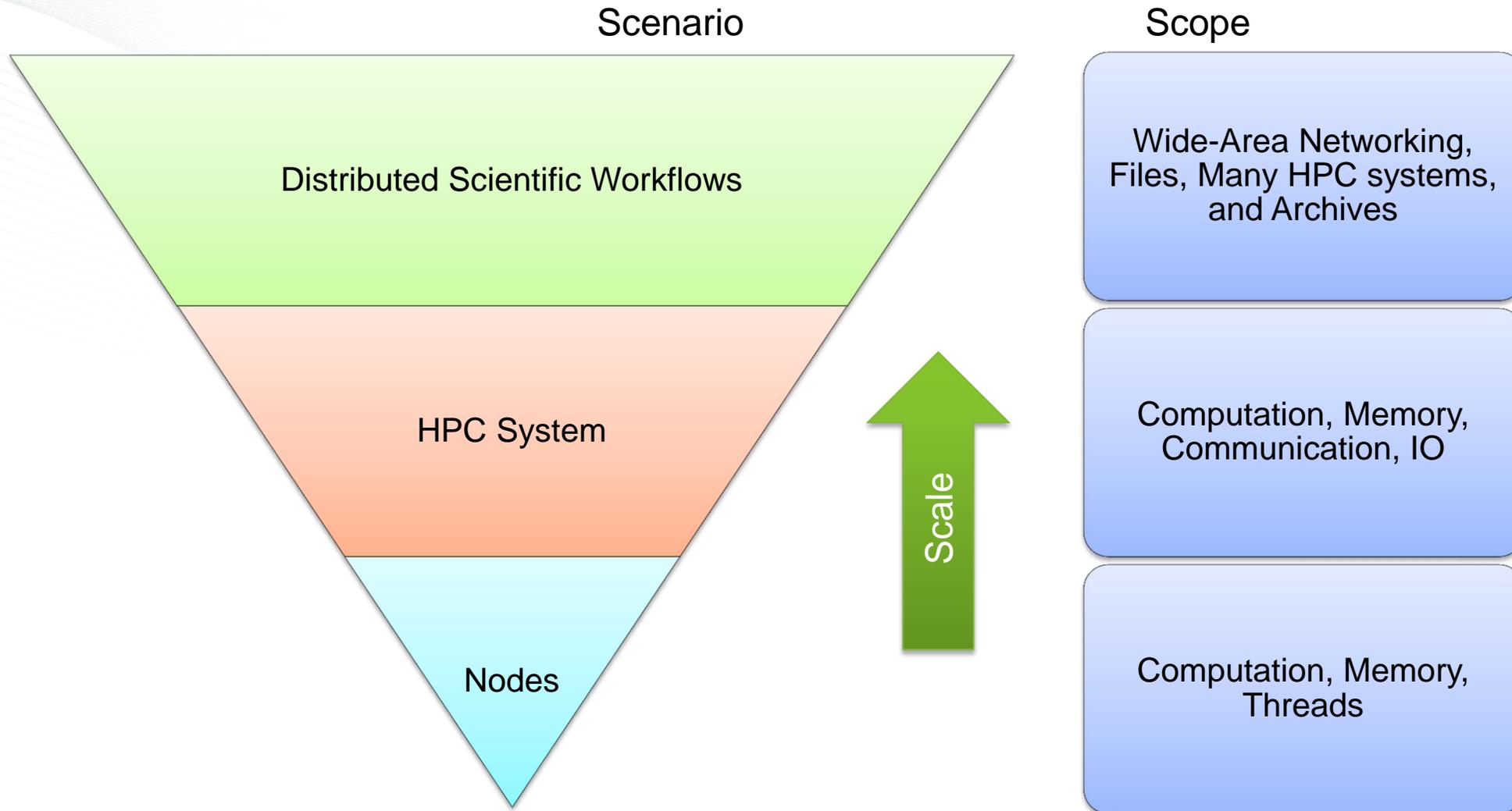
```
branch: master | aspen / models / lulesh / lulesh.aspen
jsmeredith on Sep 20, 2013 adding models
1 contributor

336 lines (288 sloc) | 9.213 kb
Raw Blame History

1 //
2 // lulesh.aspen
3 //
4 // An ASPEN application model for the LULESH 1.01 challenge problem. Based
5 // on the CUDA version of the source code found at:
6 // https://computation.llnl.gov/casc/ShockHydro/
7 //
8 param nTimeSteps = 1495
9
10 // Information about domain
11 param edgeElems = 45
12 param edgeNodes = edgeElems + 1
13
14 param numElems = edgeElems^3
15 param numNodes = edgeNodes^3
16
17 // Double precision
18 param wordSize = 8
19
20 // Element data
21 data mNodeList as Array(numElems, wordSize)
22 data mMatElemList as Array(numElems, wordSize)
23 data mNodeList as Array(8 * numElems, wordSize) // 8 nodes per element
24 data mIxi as Array(numElems, wordSize)
25 data mIxi as Array(numElems, wordSize)
26 data mIxi as Array(numElems, wordSize)
27 data mIxi as Array(numElems, wordSize)
28 data mIxi as Array(numElems, wordSize)
29 data mIxi as Array(numElems, wordSize)
30 data mIxi as Array(numElems, wordSize)
31 data mE as Array(numElems, wordSize)
32 data mP as Array(numElems, wordSize)
```

```
147 kernel CalcMonotonicQGradients {
148     execute [numElems]
149     {
150         loads [8 * indexWordSize] from nodelist
151         // Load and cache position and velocity.
152         loads/caching [8 * wordSize] from x
153         loads/caching [8 * wordSize] from y
154         loads/caching [8 * wordSize] from z
155
156         loads/caching [8 * wordSize] from xv
157         loads/caching [8 * wordSize] from yv
158         loads/caching [8 * wordSize] from zv
159
160         loads [wordSize] from v0
161         loads [wordSize] from vnew
162         // dx, dy, etc.
163         flops [90] as dp, simd
164         // delvk delxk
165         flops [9 + 8 + 3 + 30 + 5] as dp, simd
166         stores [wordSize] to delv_xeta
167         // delxi delvi
168         flops [9 + 8 + 3 + 30 + 5] as dp, simd
169         stores [wordSize] to delx_xi
170         // delxj and delvj
171         flops [9 + 8 + 3 + 30 + 5] as dp, simd
172         stores [wordSize] to delv_eta
173     }
174 }
```

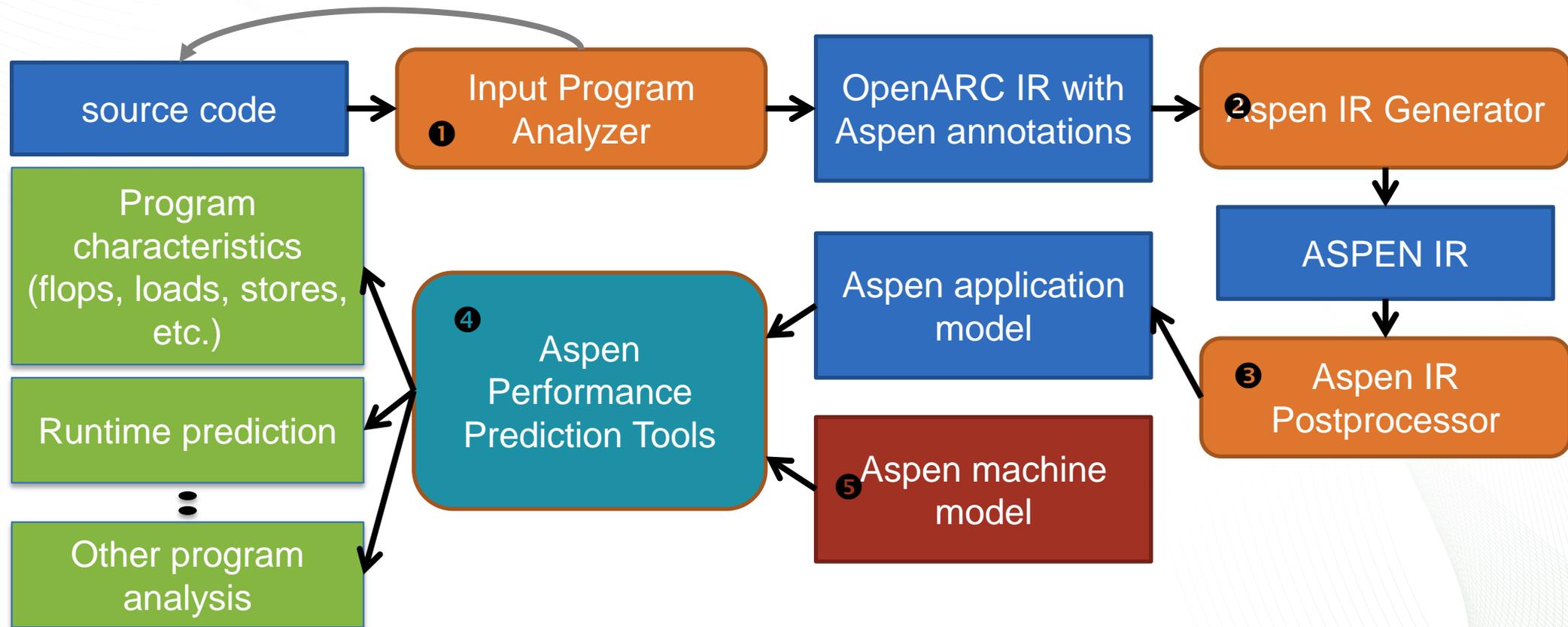
Aspen allows Multiresolution Modeling



COMPASS System Overview

- Detailed Workflow of the COMPASS Modeling Framework

Optional feedback for advanced users



MM example generated from COMPASS

```
1 int N = 1024;
2 void matmul(float *a, float *b, float *c){ int i, j, k ;
3 #pragma acc kernels loop gang copyout(a[0:(N*N)]) \
4 copyin(b[0:(N*N)],c[0:(N*N)])
5   for (i=0; i<N; i++){
6     #pragma acc loop worker
7       for (j=0; j<N; j++) { float sum = 0.0 ;
8         for (k=0; k<N; k++) {sum+=b[i*N+k]*c[k*N+j];}
9         a[i*N+j] = sum; }
10    } //end of i loop
11  } //end of matmul()
12 int main() {
13   int i; float *A = (float*) malloc(N*N*sizeof(float));
14   float *B = (float*) malloc(N*N*sizeof(float));
15   float *C = (float*) malloc(N*N*sizeof(float));
16   for (i = 0; i < N*N; i++)
17     { A[i] = 0.0F; B[i] = (float) i; C[i] = 1.0F; }
18 #pragma aspen modelregion label(MM)
19   matmul(A,B,C);
20   free(A); free(B); free(C); return 0;
21 } //end of main()
```

```
1 model MM {
2   param floatS = 4; param N = 1024
3   data A as Array((N*N), floatS)
4   data B as Array((N*N), floatS)
5   data C as Array((N*N), floatS)
6   kernel matmul {
7     execute matmul2_intracommIN
8     { intracomm [floatS*(N*N)] to C as copyin
9       intracomm [floatS*(N*N)] to B as copyin }
10  map matmul2 [N] {
11    map matmul3 [N] {
12      iterate [N] {
13        execute matmul5
14        { loads [floatS] from B as stride(1)
15          loads [floatS] from C; flops [2] as sp, simd }
16      } //end of iterate
17      execute matmul6 { stores [floatS] to A as stride(1) }
18    } // end of map matmul3
19  } //end of map matmul2
20  execute matmul2_intracommOUT
21  { intracomm [floatS*(N*N)] to A as copyout }
22  } //end of kernel matmul
23  kernel main { matmul() }
24 } //end of model MM
```

Example: LULESH (10% of 1 kernel)

```
kernel IntegrateStressForElems
{
  execute [numElem_CalcVolumeForceForElems]
  {
    loads [((1*aspen_param_int)*8)] from elemNodes as stride(1)
    loads [((1*aspen_param_double)*8)] from m_x
    loads [((1*aspen_param_double)*8)] from m_y
    loads [((1*aspen_param_double)*8)] from m_z
    loads [(1*aspen_param_double)] from determ as stride(1)
    flops [8] as dp, simd
    flops [3] as dp, simd
    stores [(1*aspen_param_double)] as stride(o)
    flops [2] as dp, simd
    stores [(1*aspen_param_double)] as stride(o)
    flops [2] as dp, simd
    stores [(1*aspen_param_double)] as stride(o)
    flops [2] as dp, simd
    loads [(1*aspen_param_double)] as stride(o)
    stores [(1*aspen_param_double)] as stride(o)
    loads [(1*aspen_param_double)] as stride(o)
    stores [(1*aspen_param_double)] as stride(o)
    loads [(1*aspen_param_double)] as stride(o)
    . . . . .
  }
}
```

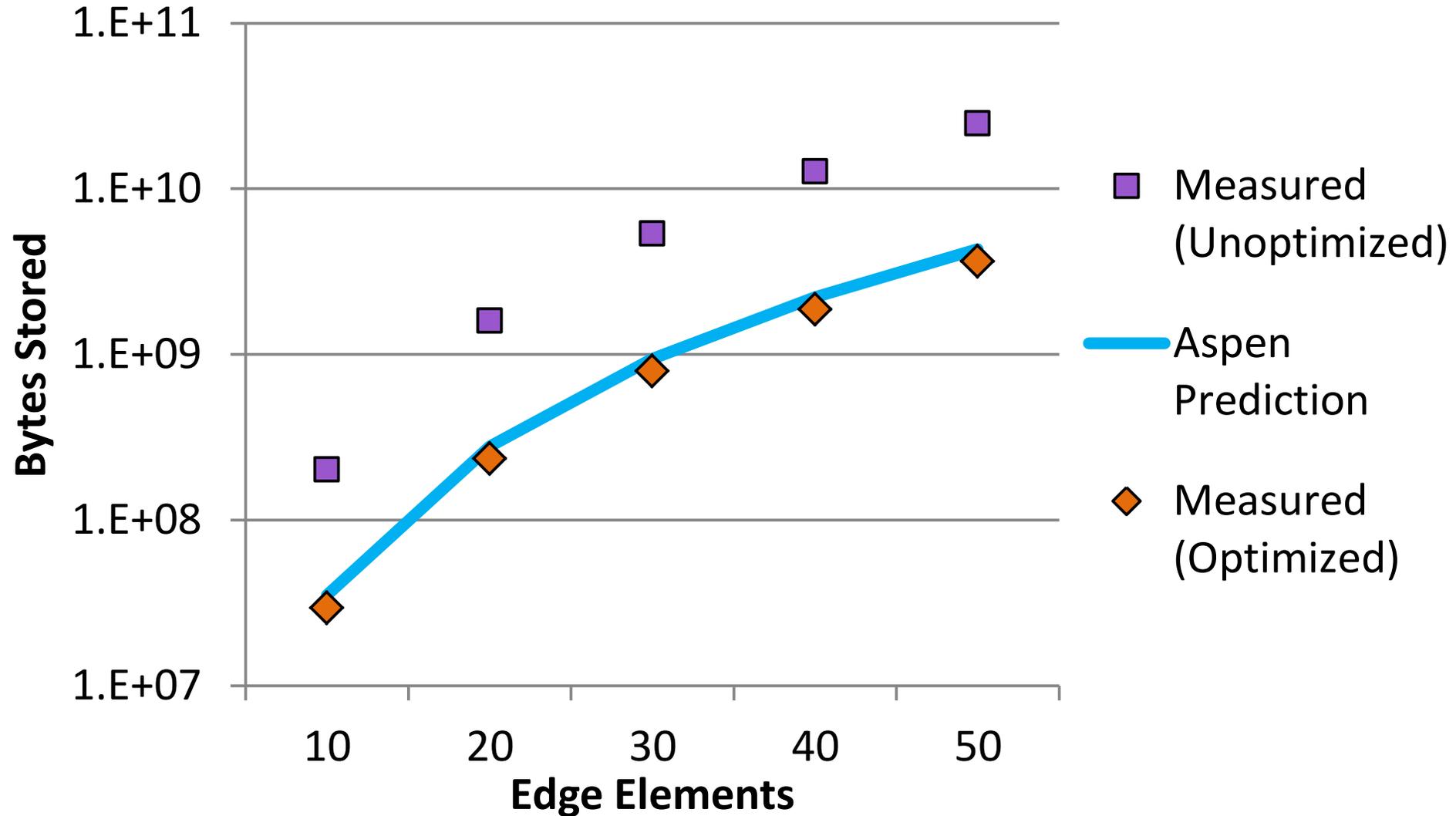
- Input LULESH program: 3700 lines of C codes
- Output Aspen model: 2300 lines of Aspen codes

Model Validation

	FLOPS	LOADS	STORES
MATMUL	15%	<1%	1%
LAPLACE2D	7%	0%	<1%
SRAD	17%	0%	0%
JACOBI	6%	<1%	<1%
KMEANS	0%	0%	8%
LUD	5%	0%	2%
BFS	<1%	11%	0%
HOTSPOT	0%	0%	0%
LULESH	0%	0%	0%

0% means that prediction fell between measurements from optimized and unoptimized runs of the code.

Model Scaling Validation (LULESH)



Example Queries

Benchmark	Runtime Order
BACKPROP	$H * O + H * I$
BFS	$nodes + edges$
CFD	$n_{elr} * n_{dim}$
CG	$n_{row} + n_{col}$
HOTSPOT	$sim_{time} * rows * cols$
JACOBI	$m_{size} * m_{size}$
KMEANS	$n_{Attr} * n_{Clusters}$
LAPLACE2D	n^2
LUD	$matrix_{dim}^3$
MATMUL	$N * M * P$
NW	max_{cols}^2
SPMUL	$size + nonzero$
SRAD	$niter * rows * cols$

Table 2: Order analysis, showing Big O runtime for each benchmark in terms of its key parameters.

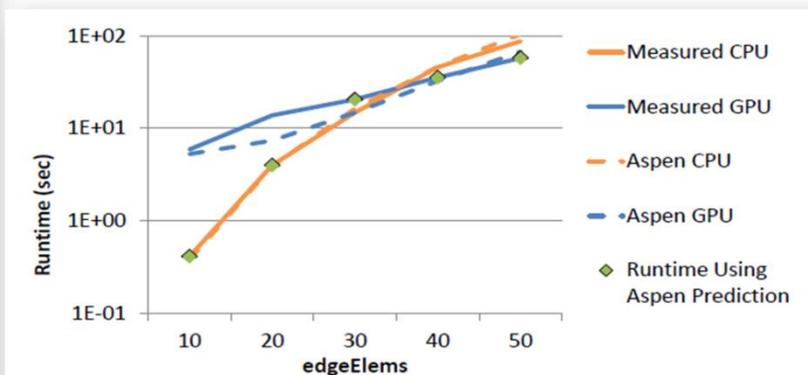
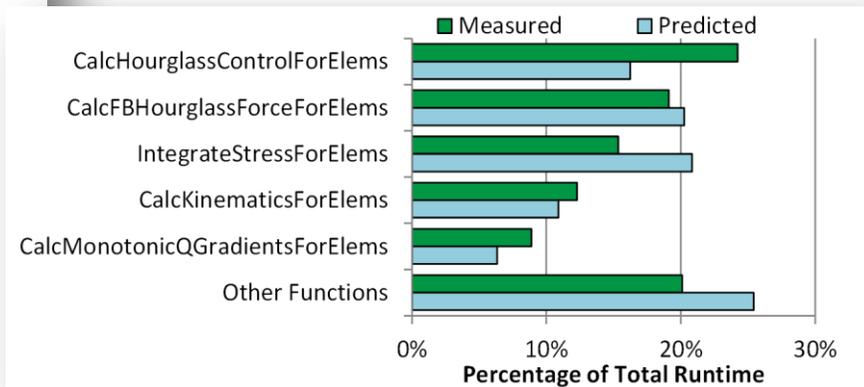


Fig. 7: Measured and predicted runtime of the entire LULESH program on CPU and GPU, including measured runtimes using the automatically predicted optimal target device at each size.

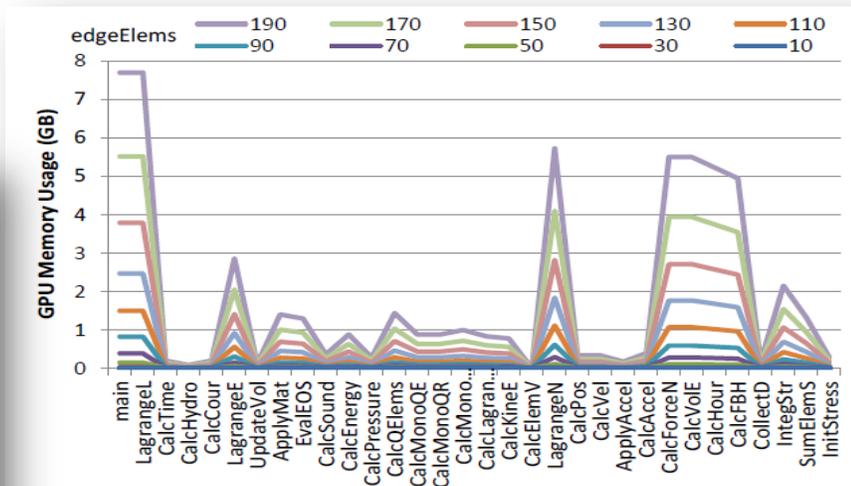


Fig. 8: GPU Memory Usage of each Function in LULESH, where the memory usage of a function is inclusive; value for a parent function includes data accessed by its child functions in the call graph.

Method Name	FLOPS/byte
InitStressTermsForElems	0.03
CalcElemShapeFunctionDerivatives	0.44
SumElemFaceNormal	0.50
CalcElemNodeNormals	0.15
SumElemStressesToNodeForces	0.06
IntegrateStressForElems	0.15
CollectDomainNodesToElemNodes	0.00
VoluDer	1.50
CalcElemVolumeDerivative	0.33
CalcElemFBHourglassForce	0.15
CalcFBHourglassForceForElems	0.17
CalcHourglassControlForElems	0.19
CalcVolumeForceForElems	0.18
CalcForceForNodes	0.18
CalcAccelerationForNodes	0.04
ApplyAccelerationBoundaryCond	0.00
CalcVelocityForNodes	0.13
CalcPositionForNodes	0.13
LagrangeNodal	0.18
AreaFace	10.25
CalcElemCharacteristicLength	0.44
CalcElemVelocityGrandient	0.13
CalcKinematicsForElems	0.24
CalcLagrangeElements	0.24
CalcMonotonicQGradientsForElems	0.46

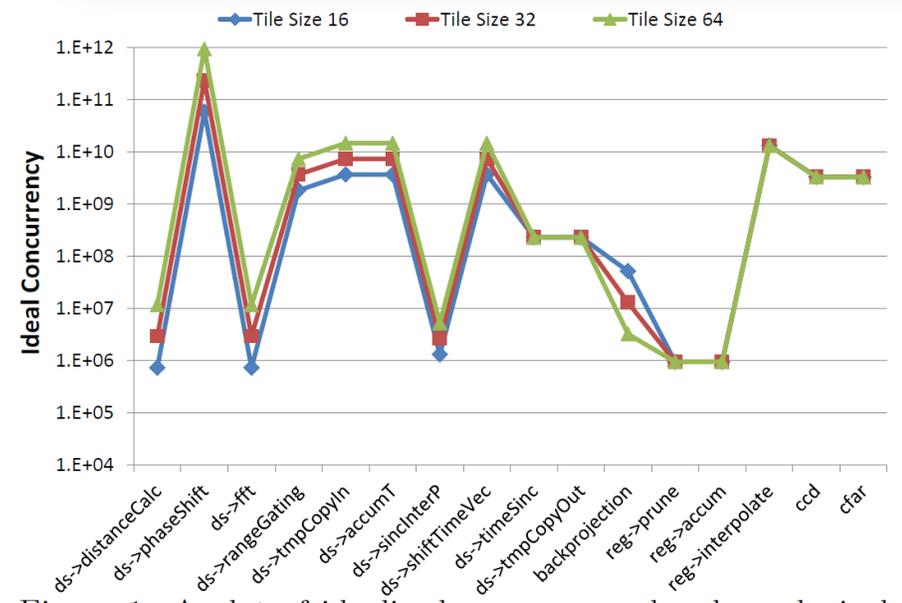
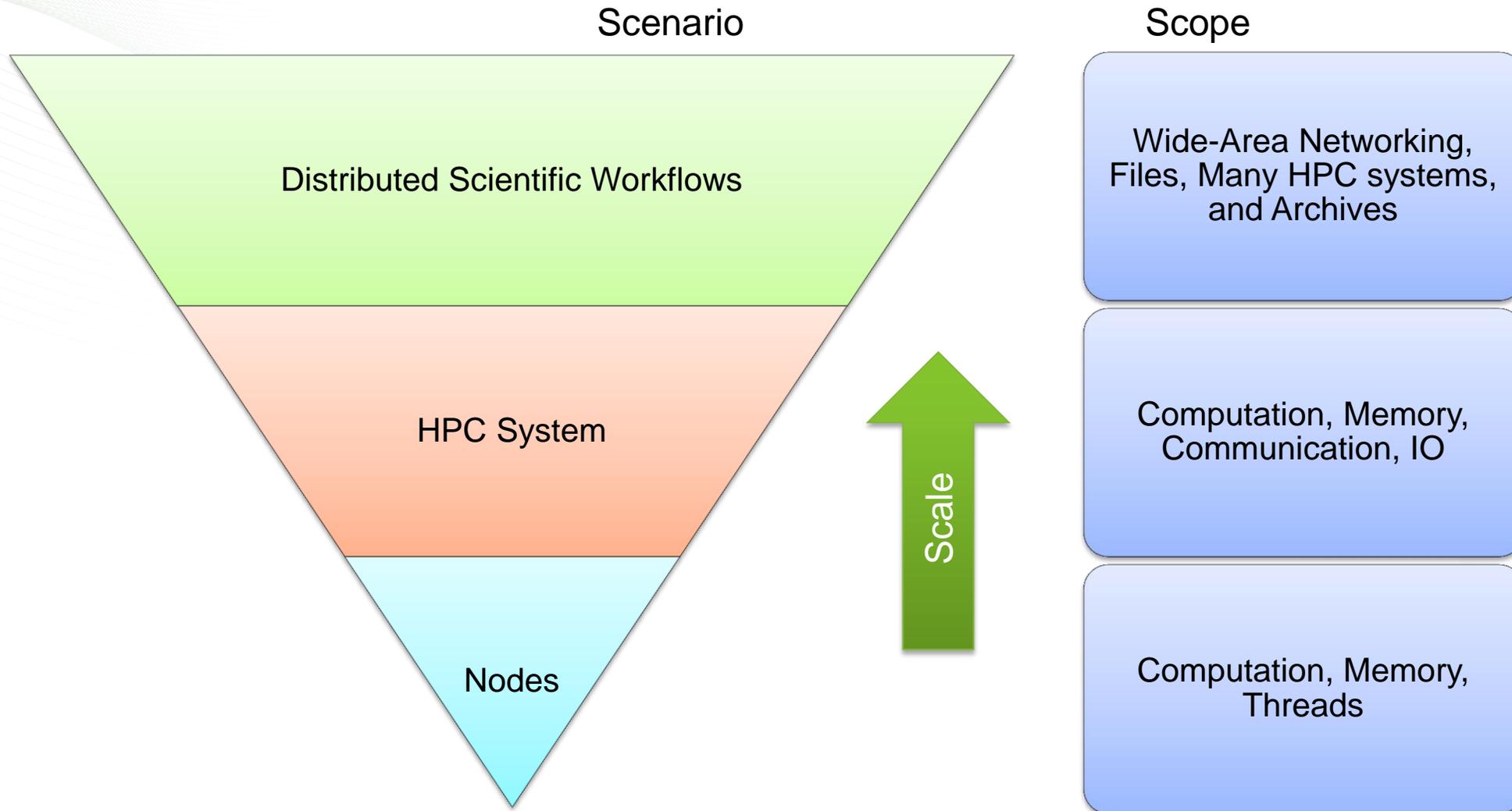
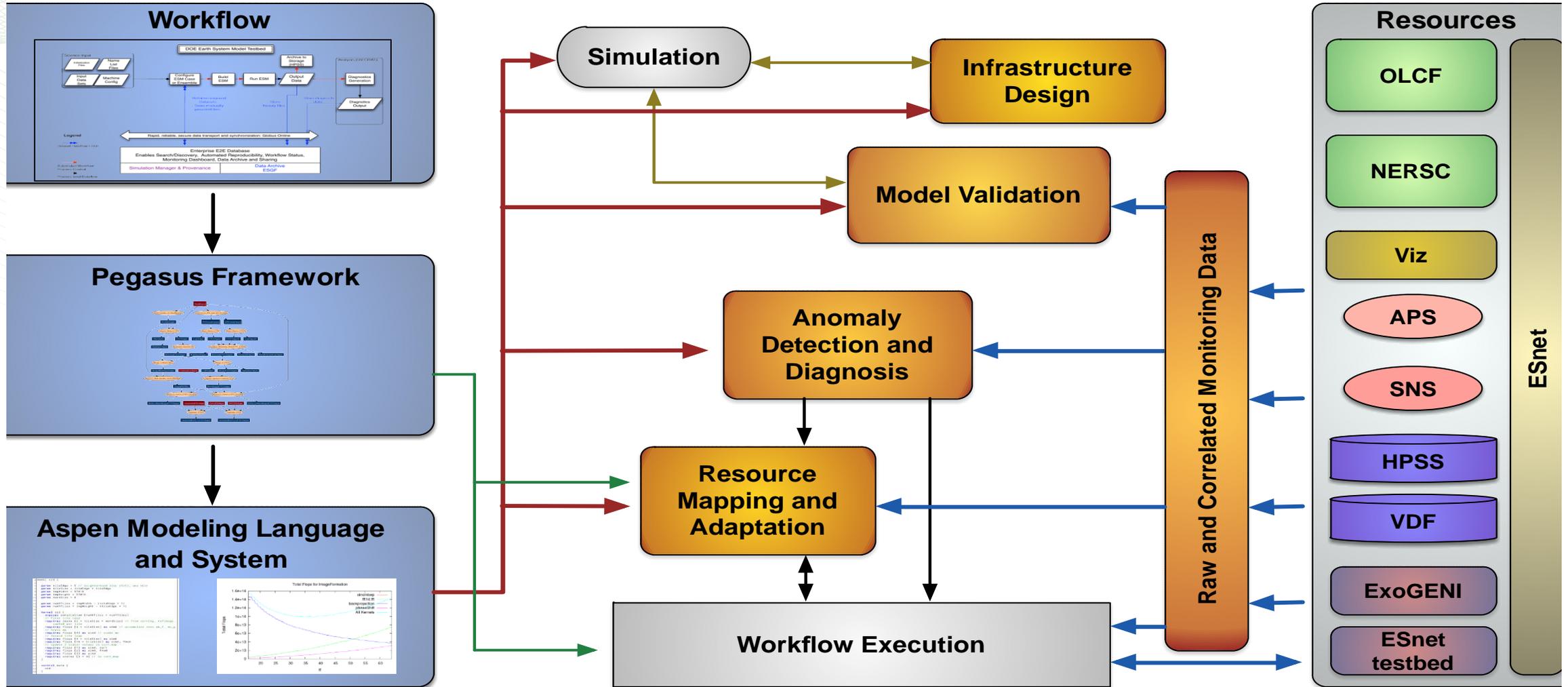


Figure 1: A plot of idealized concurrency by chronological phase in the digital spotlighting application model.

Aspen allows Multiresolution Modeling



PANORAMA Overview



Workflow
ACME
Climate
Modelin

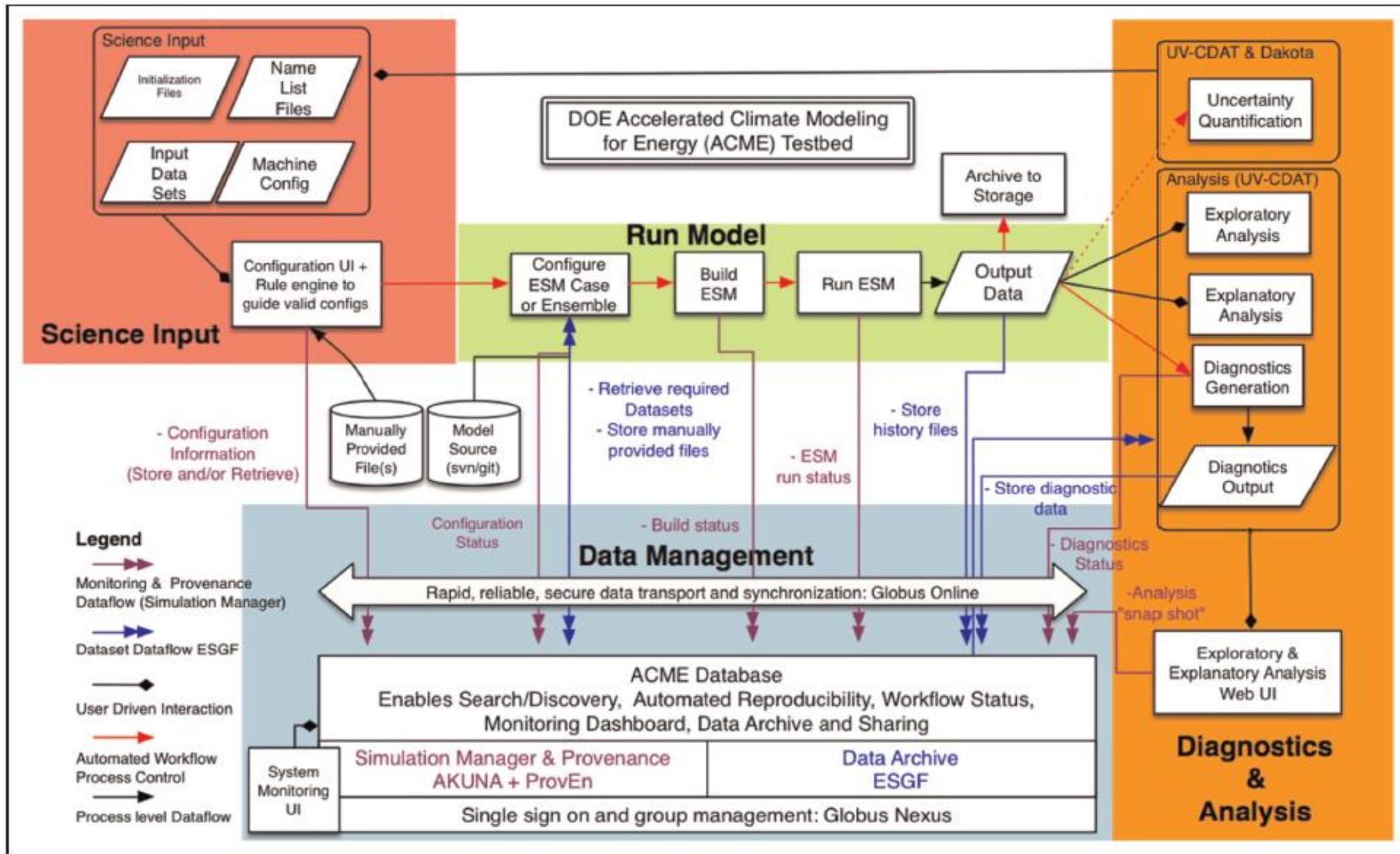
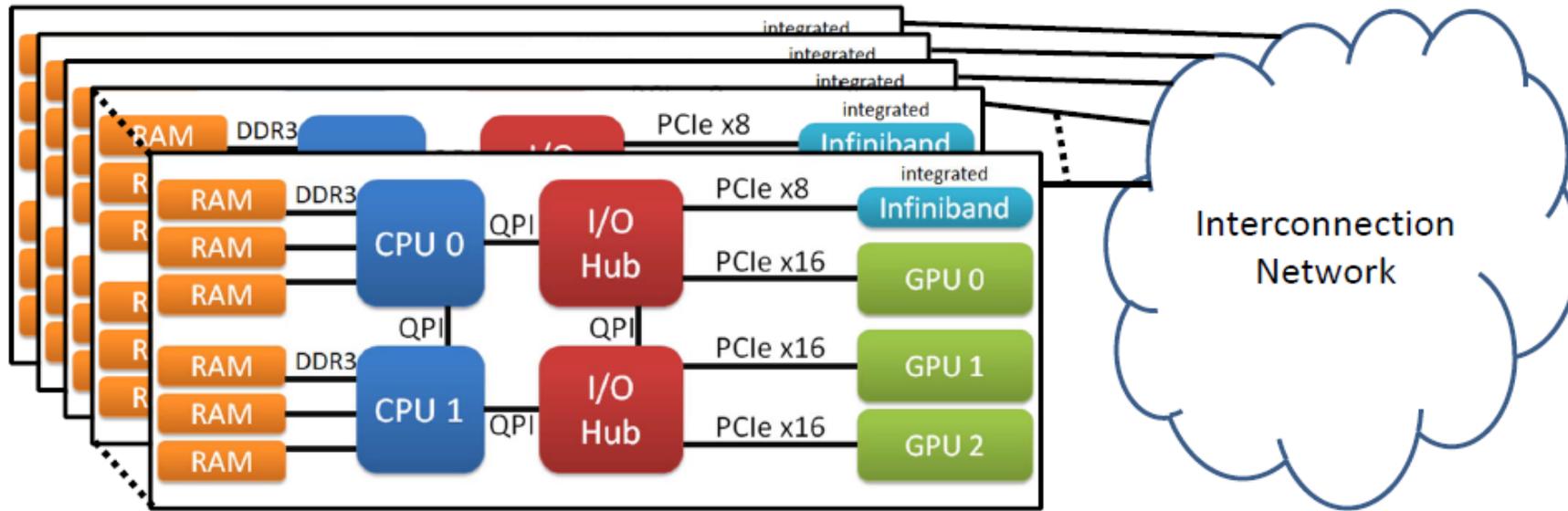


Figure 3: The complete Accelerated Climate Modeling for Energy (ACME) includes many interacting components distributed across DOE labs.

Enabling Performance Portability

Today's programming model



MPI

Low overhead

Resource contention

Locality

OpenMP, Pthreads

SIMD

NUMA

OpenACC, CUDA, OpenCL, OpenMP4, ...

Memory use,
coalescing

Data orchestration

Fine grained
parallelism

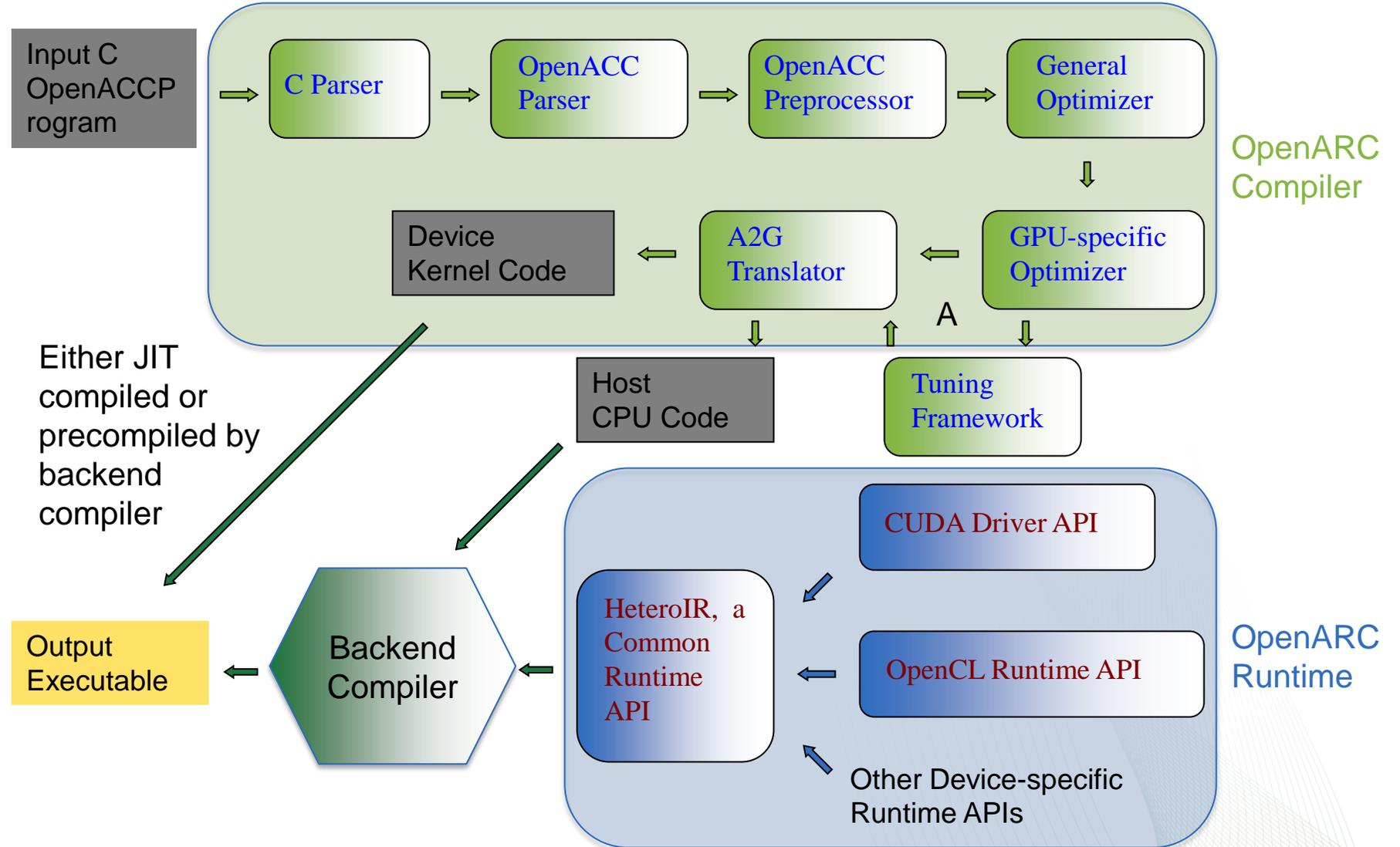
Hardware features

Contemporary Heterogeneous Architectures

Property	CUDA	GCN	MIC
Programming models	CUDA, OpenCL	OpenCL, C++ AMP	OpenCL, Cilk, TBB, LEO, OpenMP
Thread Scheduling	Hardware	Hardware	Software
Programmer Managed Cache	Yes	Yes	No
Global Synchronization	No	No	Yes
L2 Cache Type	Shared	Private per core	Private per core
L2 Total Size	Up to 1.5MB	Up to 0.5 MB	25MB
L2 Line-size	128	64	64
L1 Data Cache	Read-only + Read-write	Read-only	Read-write
Native Mode	No	No	Yes

OpenARC System Architecture

- OpenARC Compiler Framework
 - OpenACC
 - OpenMP4 Offload
- Provide common runtime APIs for various back-ends
 - CUDA
 - OpenCL
 - LLVM IR
- Open-Sourced
- High-Level Intermediate Representation (HIR)-Based
- Extensible

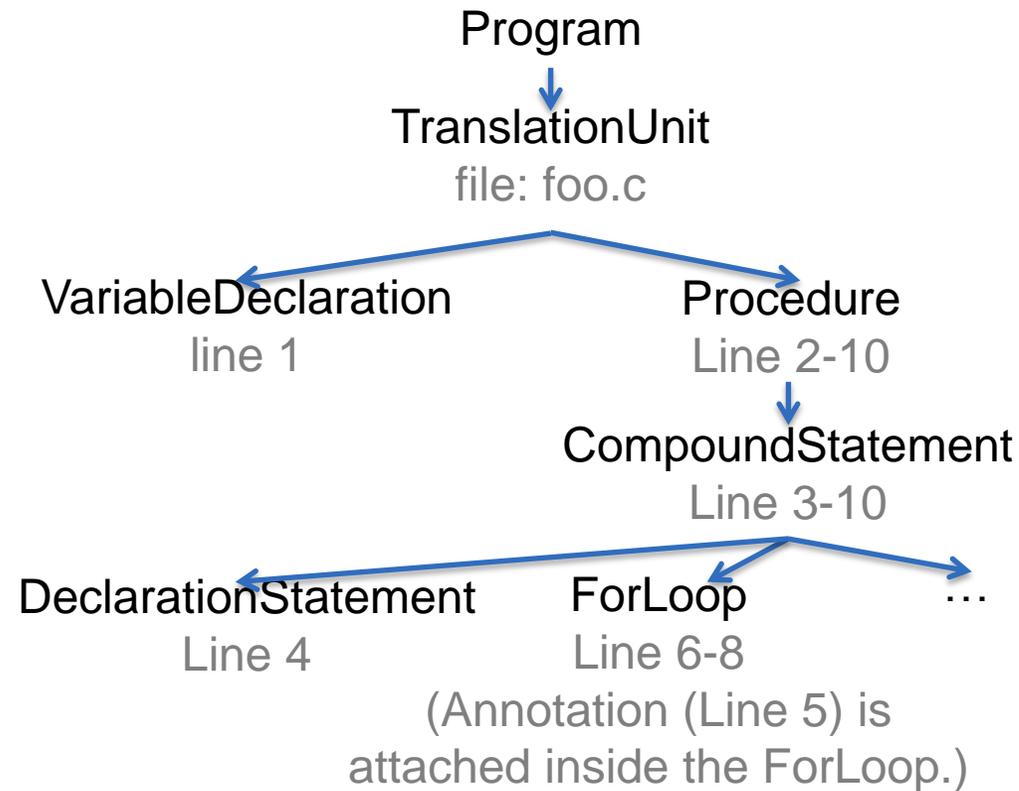


OpenARC High-Level Representation Example

- Input program

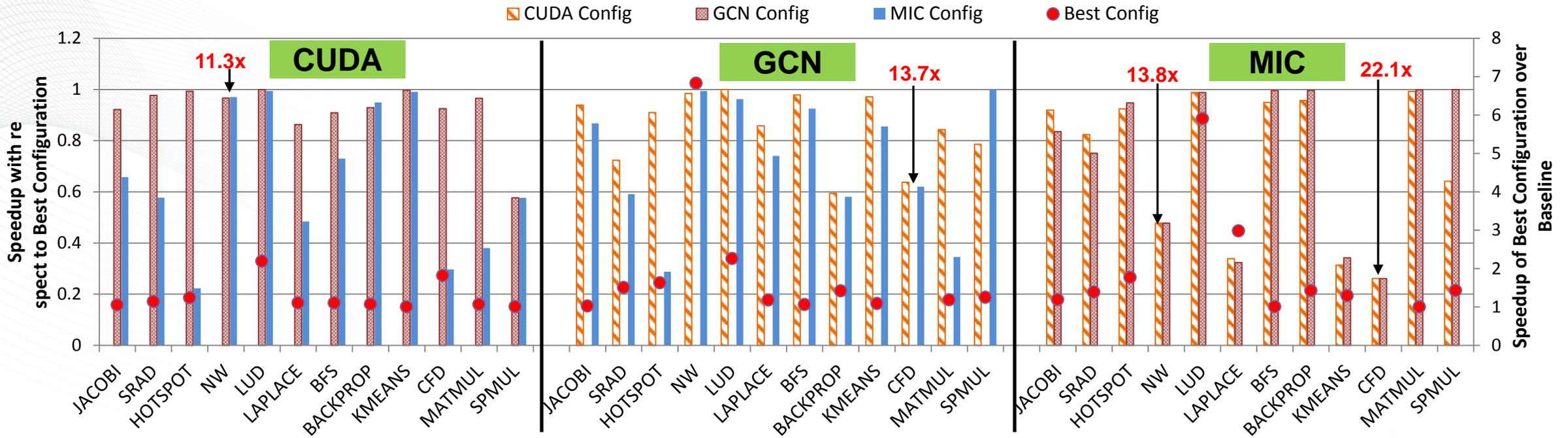
```
0 /* file: foo.c */
1 int k = 2;
2 int main(void)
3 {
4     int i, a[1024], b[1024];
5     #pragma acc kernels
6     loop copyin(b) copyout(a)
7     for( i=0; i<1024; i++ ) {
8         a[i] = k*b[i];
9     }
10 }
```

- OpenARC IR



- Common Optimizations
 - Data transfer optimizations
 - Parallel loop swap
 - Tree-based reduction generation
 - Obtaining aligned memory access
 - Loop unrolling
- Architecture specific Optimizations
 - Texture memory loading → CUDA GPUs
 - Automatic shared memory loading → CUDA, GCN GPUs
 - Pitched memory allocation → CUDA GPUs

Overall Performance Portability



- Better perf. portability among GPUs
- Lesser across GPUs and MIC
- Main reasons
 - Parallelism arrangement
 - Compiler optimizations : e.g. device-specific memories, unrolling etc.

Performance Portability Matrix

		Executed on		
		CUDA	GCN	MIC
Best Program version of	CUDA	100	84	65
	GCN	91	100	67
	MIC	58	68	100

A. Sabne, P. Sakhnagool et al., "Evaluating Performance Portability of OpenACC," in 27th International Workshop on Languages and Compiler for Parallel Computing (LCPC) Portland, Oregon, 2014

Automating selection of optimizations based on machine model

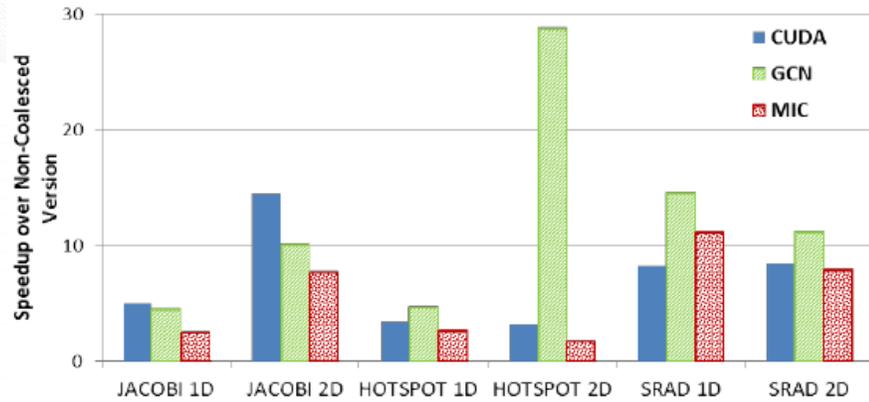


Figure 5: Memory Coalescing Benefits on Different Architectures : MIC is impacted the least by the non-coalesced accesses

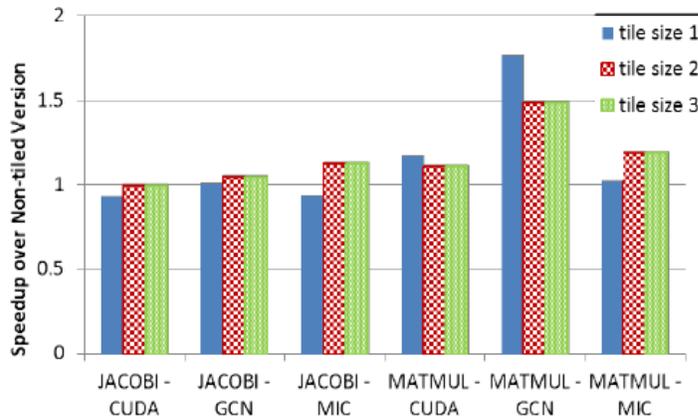


Figure 7: Impact of Tiling Transformation : *MATMUL* shows higher benefits than *JACOBI* owing to more contiguous accesses

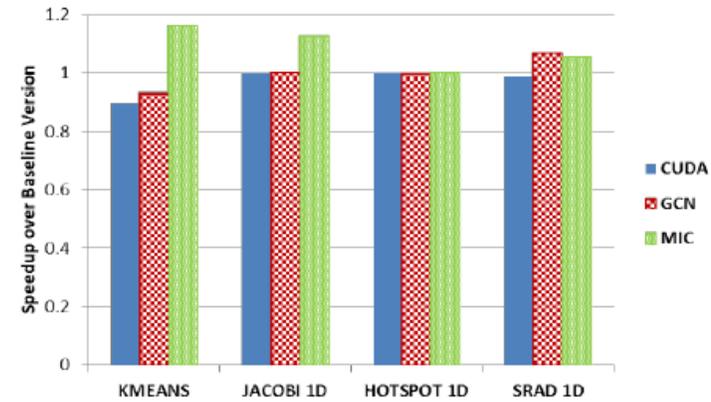


Figure 9: Effects of Loop Unrolling - MIC shows benefits on unrolling

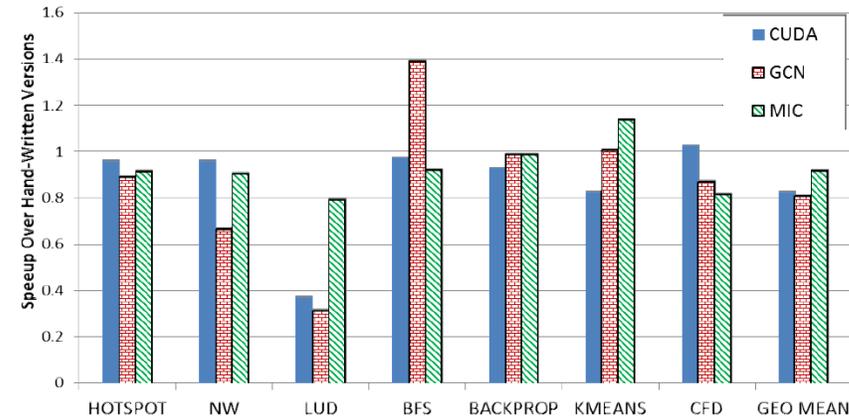


Fig. 11: Comparison of hand-written CUDA/OpenCL programs against auto-tuned OpenARC code versions : Tuned OpenACC programs perform reasonably well against hand-written codes

Overview

- Our community has major challenges in HPC as we move to extreme scale
 - Power, Performance, Resilience, Productivity
 - New technologies emerging to address some of these challenges
 - Heterogeneous computing
 - Multimode memory systems including nonvolatile memory
 - Not just HPC: Most uncertainty in at least two decades
 - Exascale includes even more diversity and uncertainty
- **We need performance prediction and portability tools now more than ever!**
- Aspen is a tool for structured design and analysis
 - Co-design applications and architectures for performance, power, resiliency
 - Automatic model generation
 - Scalable to distributed scientific workflows
- OpenARC research compiler is a vehicle for
 - Understanding how to automate platform specific optimizations
 - Developing performance portable code

Acknowledgments

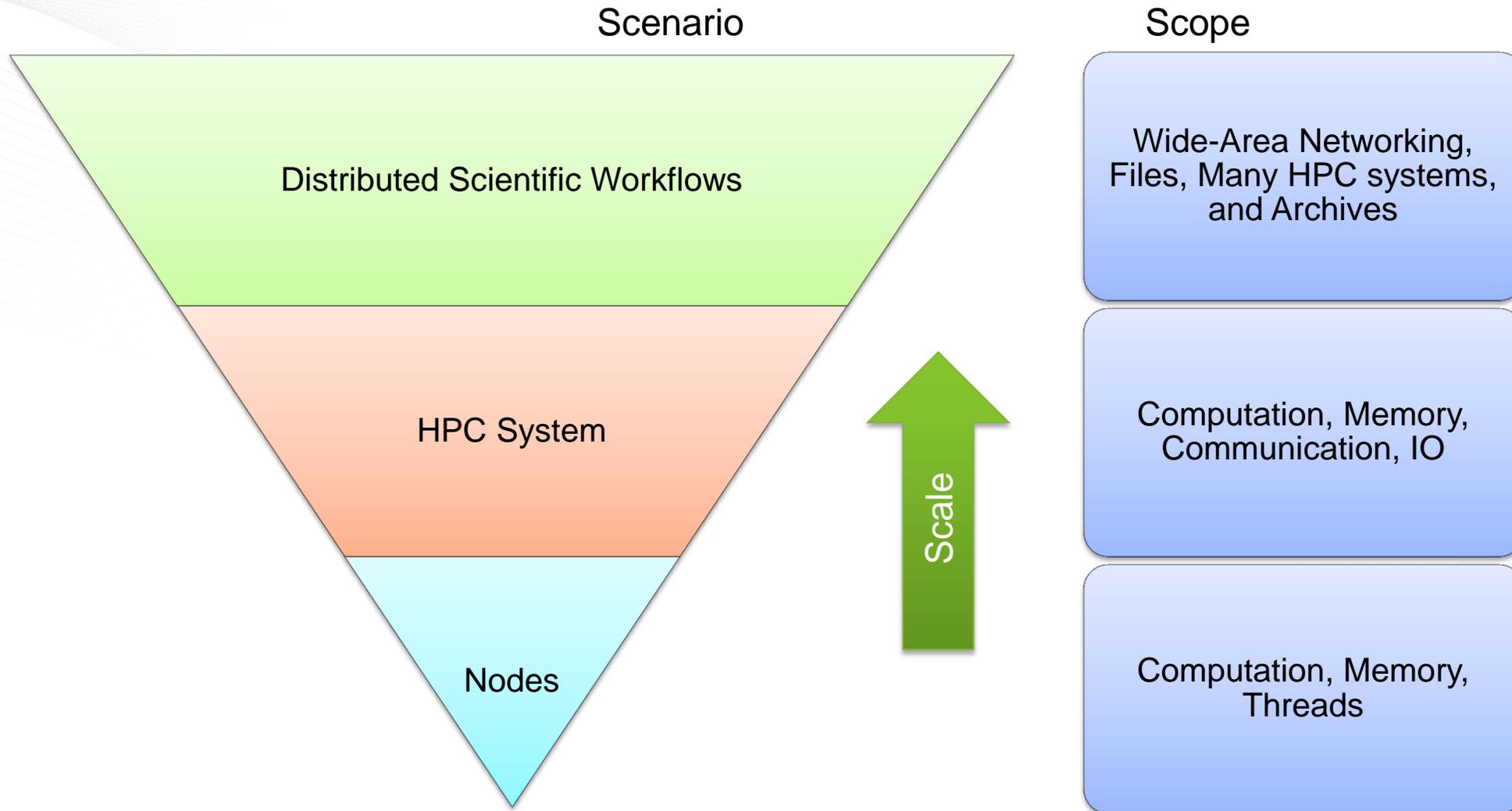
- Contributors and Sponsors

- Future Technologies Group: <http://ft.ornl.gov>
- US Department of Energy Office of Science
 - DOE Vancouver Project: <https://ft.ornl.gov/trac/vancouver>
 - DOE Blackcomb Project: <https://ft.ornl.gov/trac/blackcomb>
 - DOE ExMatEx Codesign Center: <http://codesign.lanl.gov>
 - DOE Cesar Codesign Center: <http://cesar.mcs.anl.gov/>
 - DOE Exascale Efforts: <http://science.energy.gov/ascr/research/computer-science/>
- Scalable Heterogeneous Computing Benchmark team: <http://bit.ly/shocmarx>
- US National Science Foundation Keeneland Project: <http://keeneland.gatech.edu>
- US DARPA
- NVIDIA CUDA Center of Excellence



Performance Modeling for Distributed Scientific Workflows

Aspen allows Multiresolution Modeling



Workflow: SNS

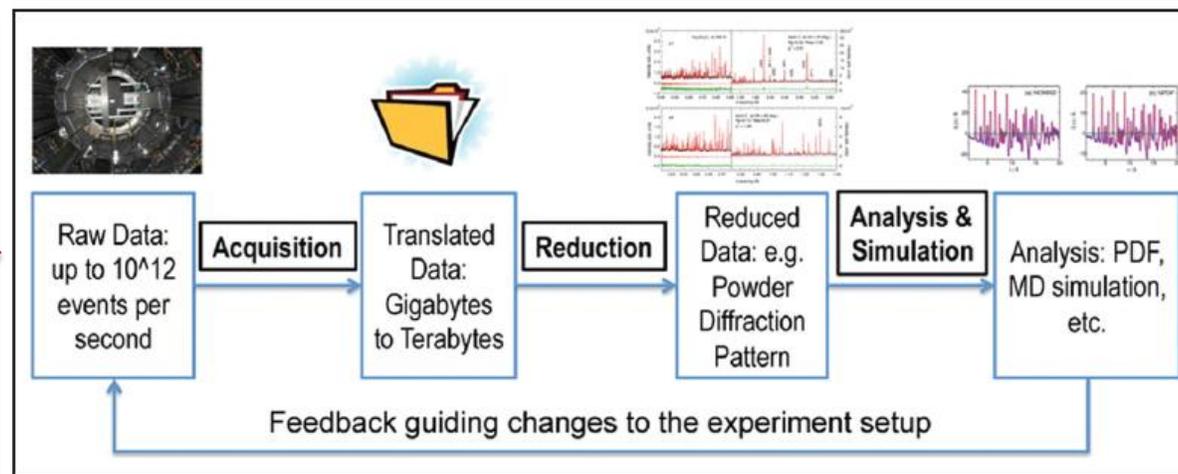
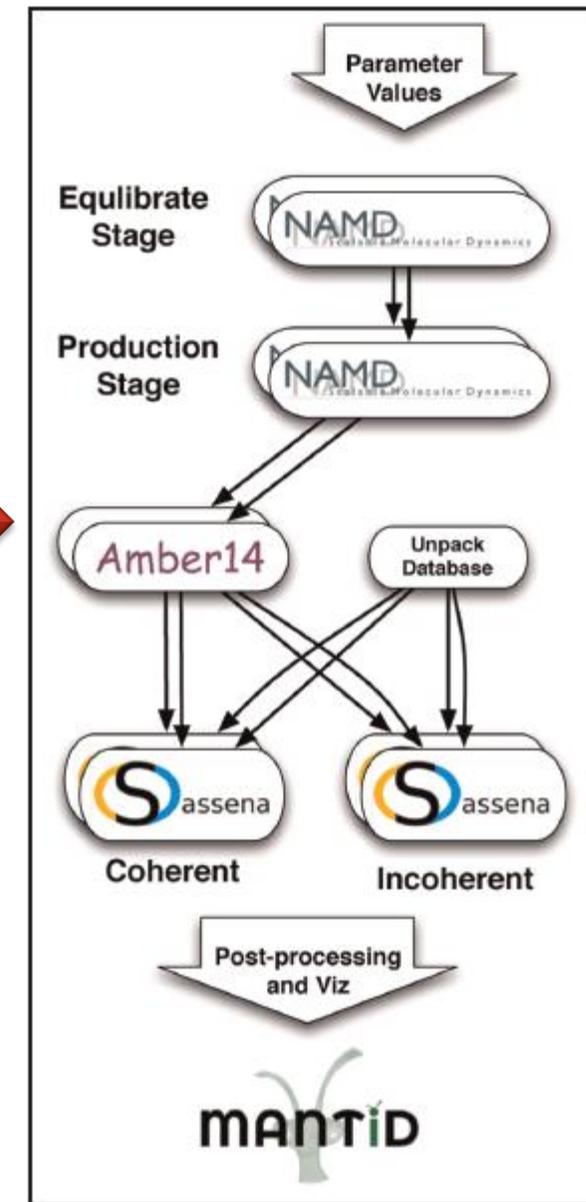
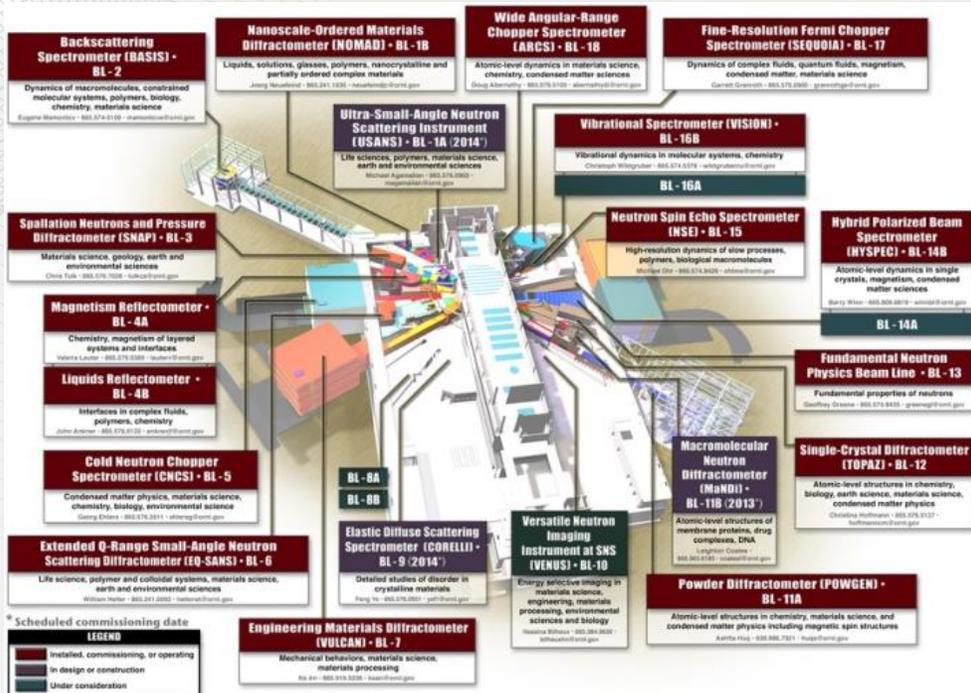
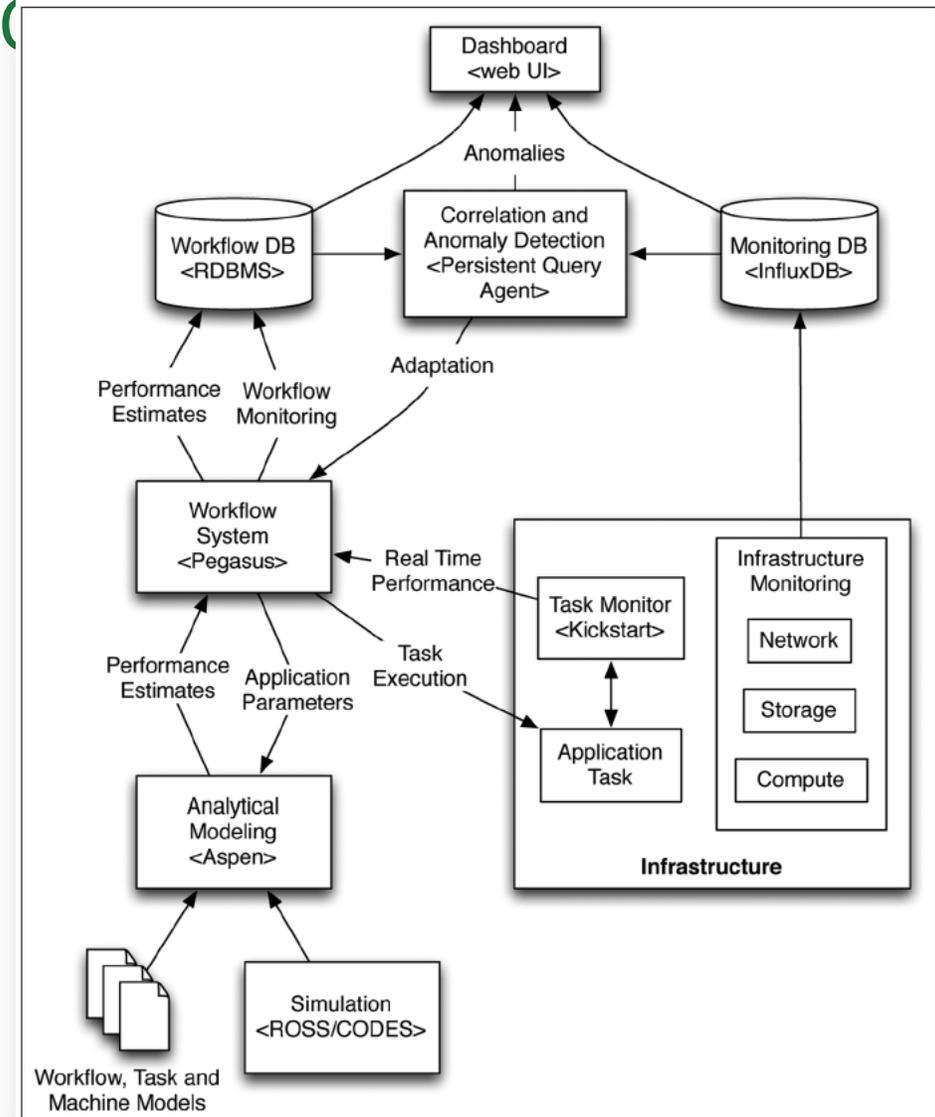


Figure 2: The SNS refinement workflow executes a parameter sweep of molecular dynamics and neutron scattering simulations to optimize the value for a target parameter to fit experimental data.

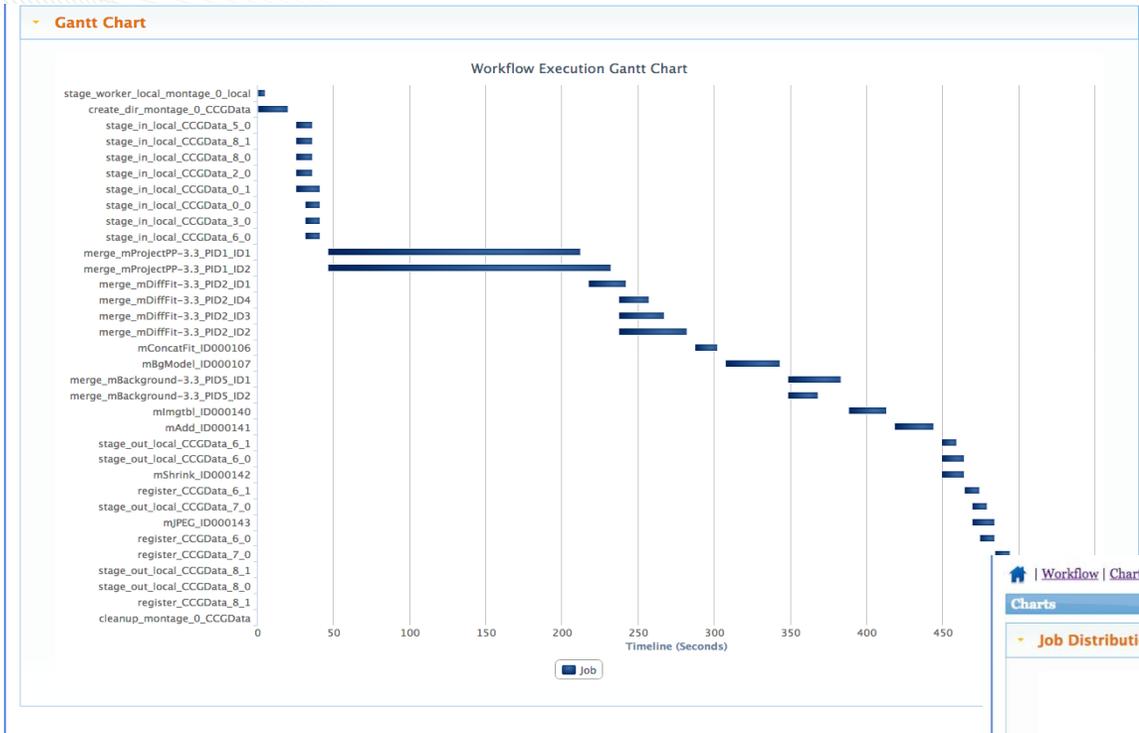
Automatically Generate Aspen from Pegasus DAX; Use Aspen Predictions to Inform/Monitor Decision

```
1 kernel main
2 {
3   par {
4     seq {
5       call namd_eq_200()
6       call namd_prod_200()
7     }
8     seq {
9       call namd_eq_290()
10      call namd_prod_290()
11    }
12  }
13  par {
14    call unpack_database()
15    call ptraj_200()
16    call ptraj_290()
17  }
18  par {
19    call sassena_incoh_200()
20    call sassena_coh_200()
21    call sassena_incoh_290()
22    call sassena_coh_290()
23  }
24 }
```

Listing 1: Automatically generated Aspen model for sample SNS workflow.



Workflow Monitoring Dashboard – *pegasus-dashboard*



Workflow | Statistics

Statistics

Workflow Wall Time	8 mins 53 secs
Workflow Cumulative Job Wall Time	1 min 59 secs
Cumulative Job Walltime as seen from Submit Side	4 mins 18 secs
Workflow Retries	0

Workflow Statistics

Job Breakdown Statistics

Show 50 entries Search: []

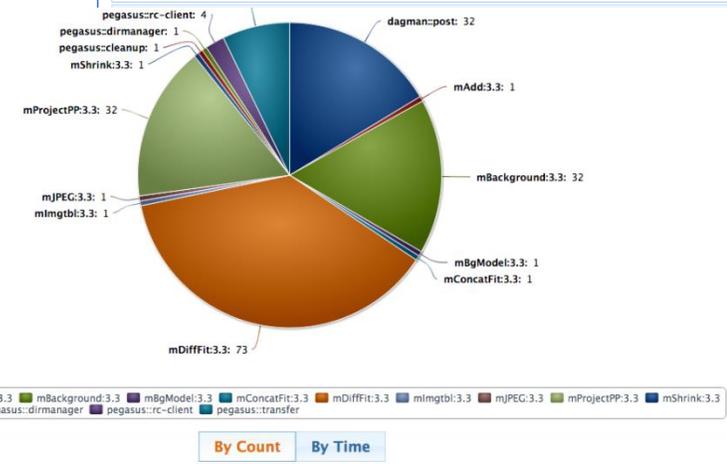
Transformation	Count	Succeeded	Failed	Min	Max	Mean	Total
dagman::post	32	32	0	5	6	5.063	162
mAdd:3.3	1	1	0	1.203	1.203	1.203	1.203
mBackground:3.3	32	32	0	0.054	0.197	0.130	4.174
mBgModel:3.3	1	1	0	18.701	18.701	18.701	18.701
mConcatFit:3.3	1	1	0	1.033	1.033	1.033	1.033
mDiffFit:3.3	73	73	0	0.048	0.226	0.103	7.492
mImgtbl:3.3	1	1	0	0.107	0.107	0.107	0.107
mJPEG:3.3	1	1	0	0.523	0.523	0.523	0.523
mProjectPP:3.3	32	32	0	0.915	0.978	0.926	29.633
mShrink:3.3	1	1	0	0.485	0.485	0.485	0.485
pegasus::cleanup	1	1	0	5	5	5	5
pegasus::dirmanager	1	1	0	10	10	10	10
pegasus::rc-client	4	4	0	0.706	0.868	0.783	3.134
pegasus::transfer	14	14	0	0	5.229	2.724	38.135

Showing 1 to 14 of 14 entries First Previous 1 Next Last

Workflow | Charts

Charts

Job Distribution



Status, statistics, timeline of jobs

Helps pinpoint errors