

Exceptional service in the national interest



Resource Management Challenges in the Era of Extreme Heterogeneity

Ron Brightwell, R&D Manager

Scalable System Software Department



Sandia National Laboratories is a multi-mission laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000. SAND NO. 2011-XXXXP

Outline

- Key takeaways
- Brief explanation of the national labs
- Extreme Heterogeneity Summit
- Extreme Heterogeneity Workshop
- Priority research directions for resource management
- List of issues and concerns

Key Takeaways

- A recent ASCR workshop on Extreme Heterogeneity has identified several key challenges and potential research directions in the following areas:
 - Programming environments
 - Software development, sustainability, and productivity
 - Operating systems and resource management
 - Data management, analytics, and workflows
 - Architecture modeling and simulation
- This talk will expand on the OS/RM challenges
- Results of the workshop are being compiled in a report which may (or may not) be used as a basis for future ASCR program investments

Funding Models at the National Labs

NNSA

- LLNL, LANL, SNL
 - Advanced Simulation and Computing (ASC)
 - Program elements
 - Integrated Codes (IC)
 - Physics and Engineering Models (P&EM)
 - Verification and Validation (V&V)
 - Facilities, Operations, and User Support (FOUS)
 - Computational Systems and Software Environment (CSSE)
 - Advanced Technology Design and Mitigation (ATDM)
 - Stockpile stewardship mission
 - Direct funding to accomplish mission



DOE

- ANL, ORNL, LBNL, PNNL, BNL, ...
 - Office of Science
 - Program Offices
 - Advanced Scientific Computing Research (ASCR)
 - Basic Energy Sciences (BES)
 - Biological and Environment Research (BER)
 - Fusion Energy Sciences (FES)
 - High Energy Physics (HEP)
 - Nuclear Physics (NP)
 - Science mission
 - Program funding model
 - Competitive proposals

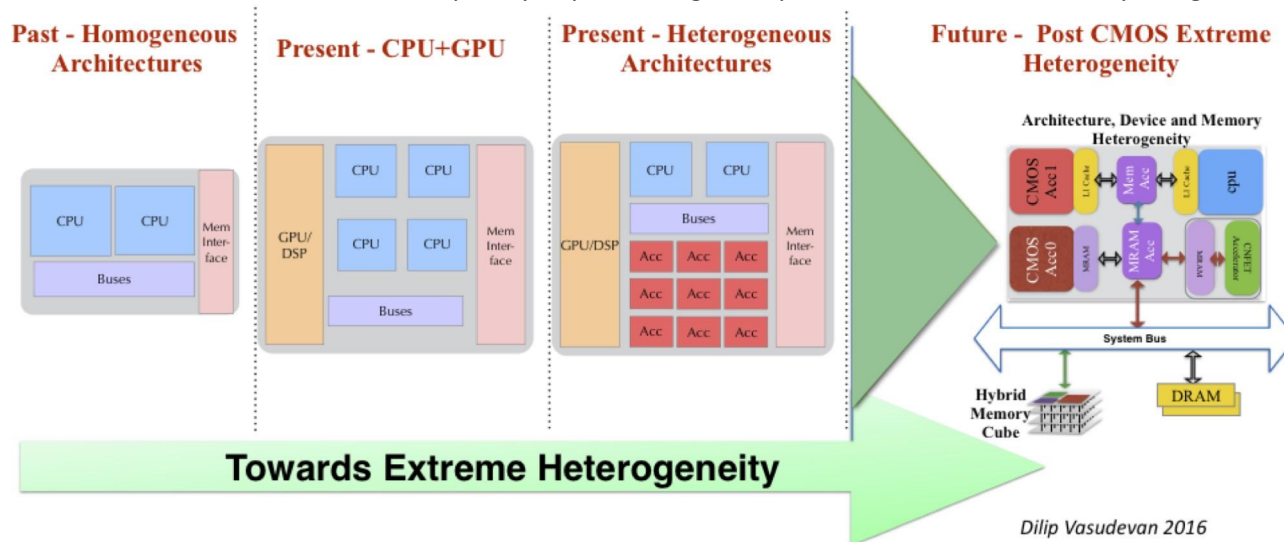


ASCR Extreme Heterogeneity Summit

- June 8-9, 2017
- Participants
 - Jeffrey Vetter (ORNL) Rob Ross (ANL), Pat McCormick (LANL), Katie Antypas (LBL), John Shalf (LBL), David Donofrio (LBL), Maya Gokhale (LLNL), Ron Brightwell (SNL), Travis Humble (ORNL), ShinJae Yoo (BNL), Catherine Schuman (ORNL)
- Purpose
 - Determine whether workshop on Extreme Heterogeneity is needed
 - If so, begin initial planning phase for workshop
- Goals
 - Come to agreement on the definition of Extreme Heterogeneity
 - Determine topics to be addressed at the workshop
 - Develop a rough agenda
 - Identify key participants
 - Write a report summarizing the Summit

The Challenge of Heterogeneity

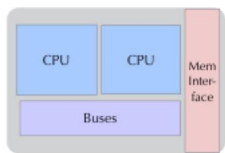
- **“A challenge of heterogeneity is how to build large systems comprised of massive numbers of these already heterogeneous systems”** Bob Colwell (former Intel chip architect and DARPA MTO Director)
- **If ASCR does not confront these challenges through new research**
 - HPC is consigned to only modest improvements beyond exascale
 - Complexity will make code maintenance impractical or unsustainable in the long term
 - Overall: cost/complexity impedes long-term pursuit of scientific discovery using HPC



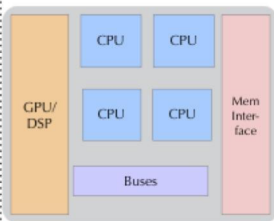
The Challenge of Heterogeneity

- **“A challenge of heterogeneity is how to build large systems comprised of massive numbers of these already heterogeneous systems”** Bob Colwell (former Intel chip architect and DARPA MTO Director)
- **If ASCR does not confront these challenges through new research**
 - HPC is consigned to only modest improvements beyond exascale
 - Complexity will make code maintenance impractical or unsustainable in the long term
 - Overall: cost/complexity impedes long-term pursuit of scientific discovery using HPC

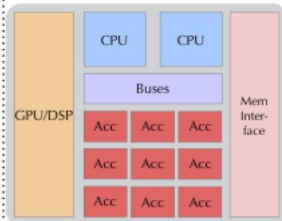
Past - Homogeneous Architectures



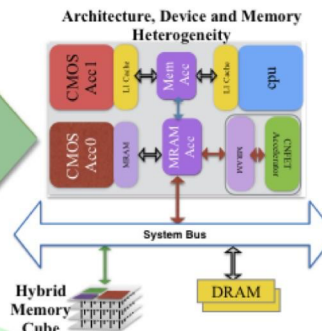
Present - CPU+GPU



Present - Heterogeneous Architectures



Future - Post CMOS Extreme Heterogeneity



Towards Extreme Heterogeneity

Past 30 years of Parallel systems
(1,000,000,000x of scaling)

Pre-Exascale
(Titan/Summit)

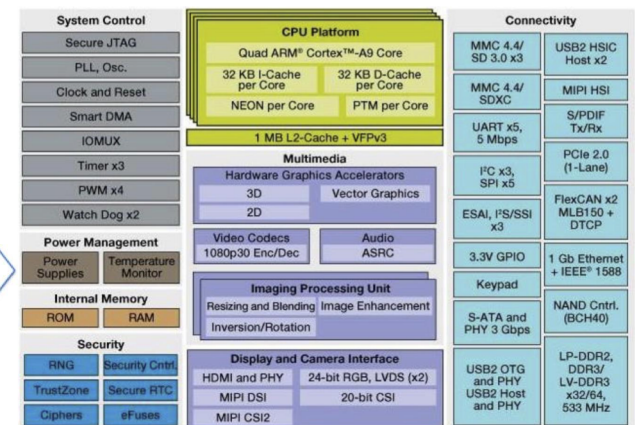
Exascale
(A21/Coral2)

Post-Exascale
(???)

Dilip Vasudevan 2016

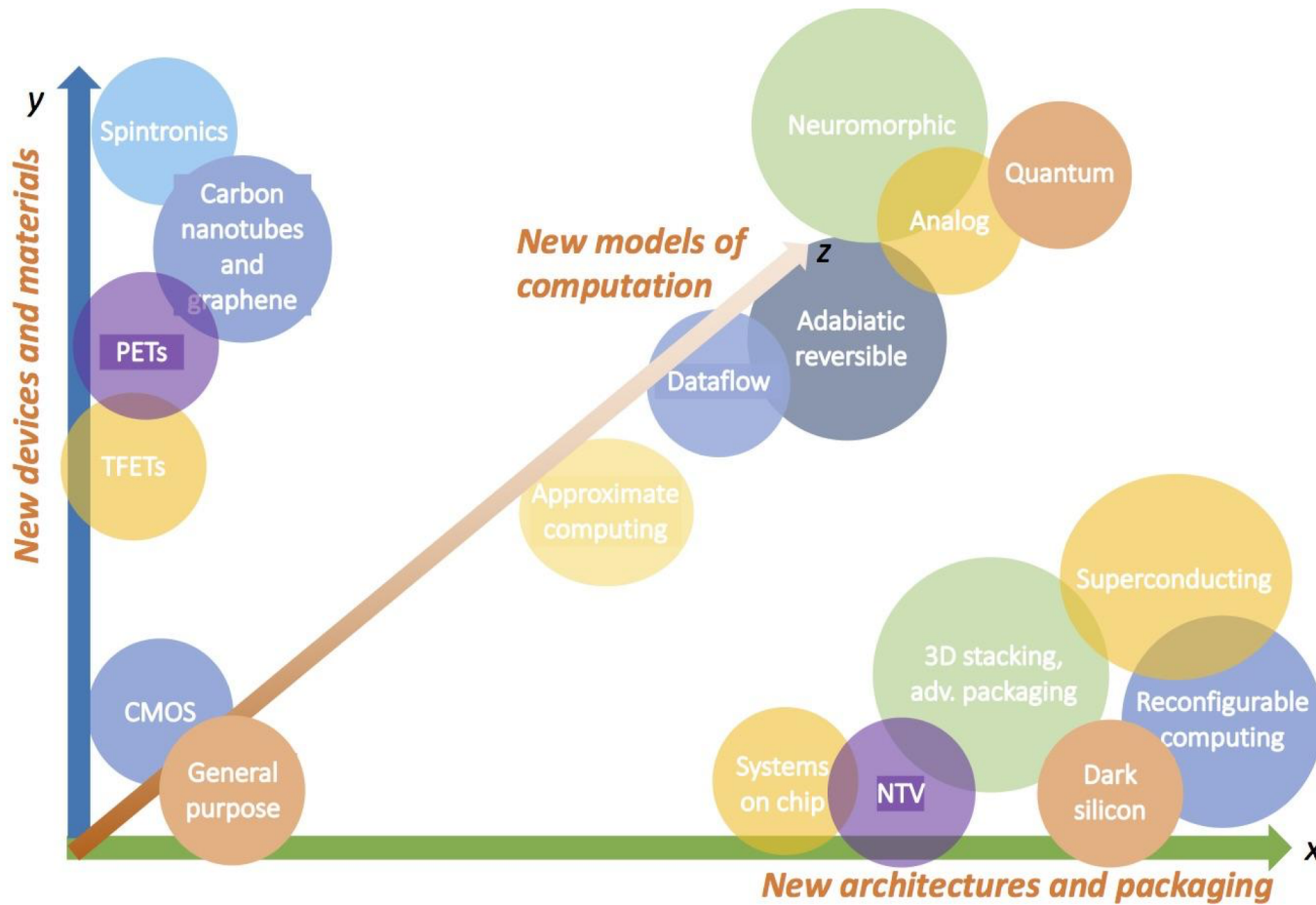
This is already happening TODAY!

Below is a SmartPhone SoC circa 2016
Dozens of kinds of integrated HW acceleration

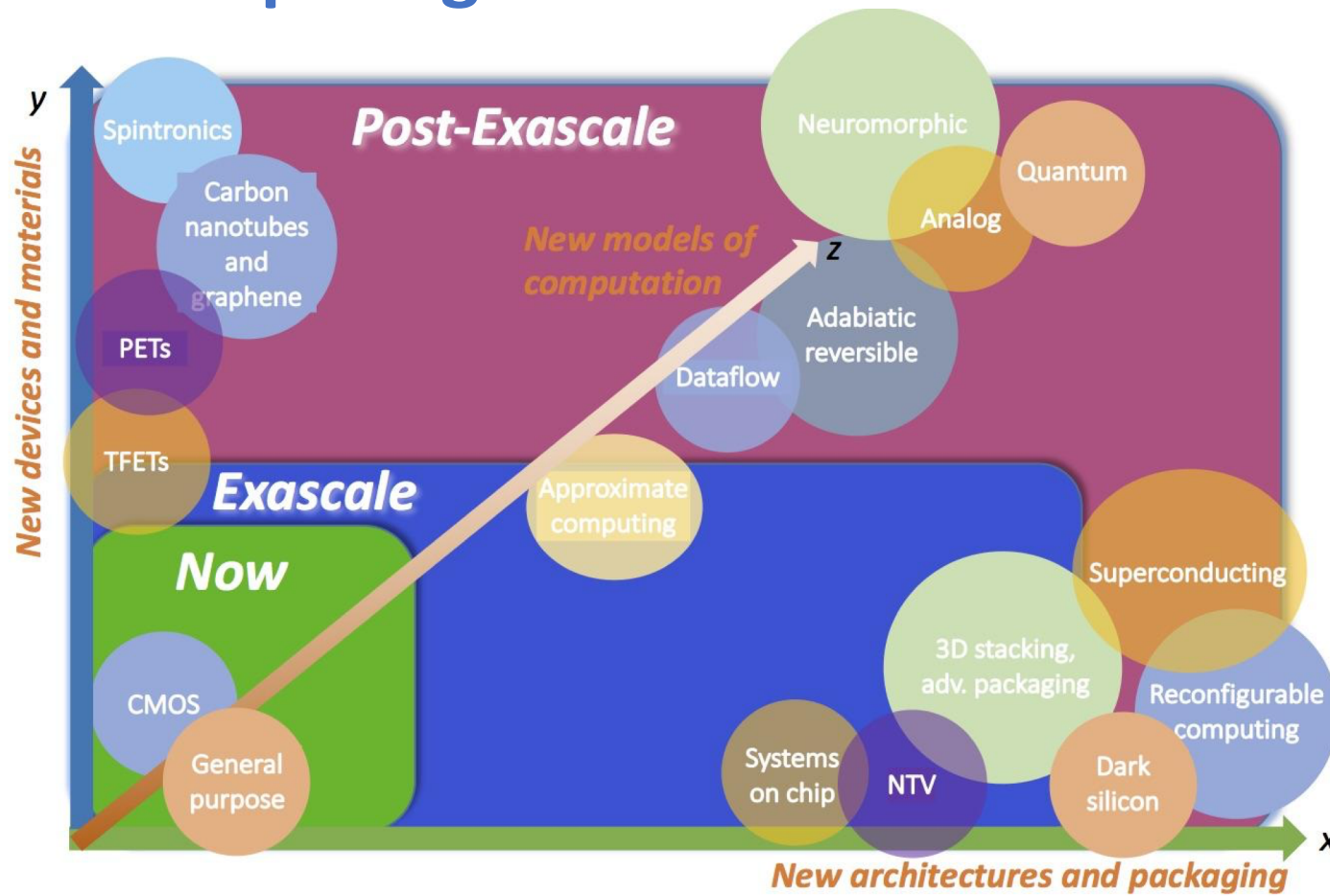


10+ years to make GPU accelerators usable for science
Will it take us 100 years to get 10 more of them usable?
Or will HPC fall behind the rest of the computing industry?

Future of Computing

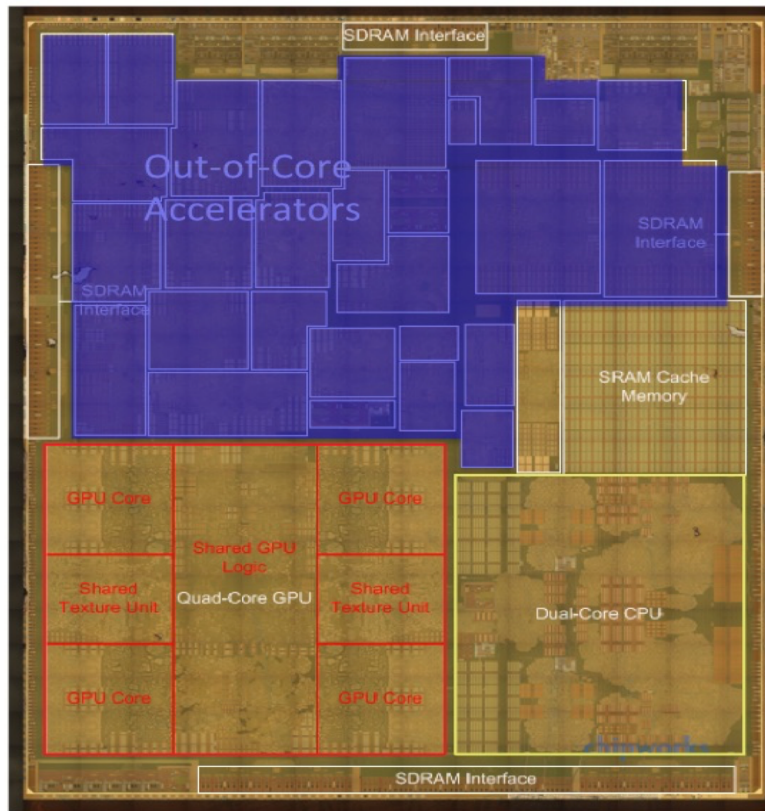


Future of Computing

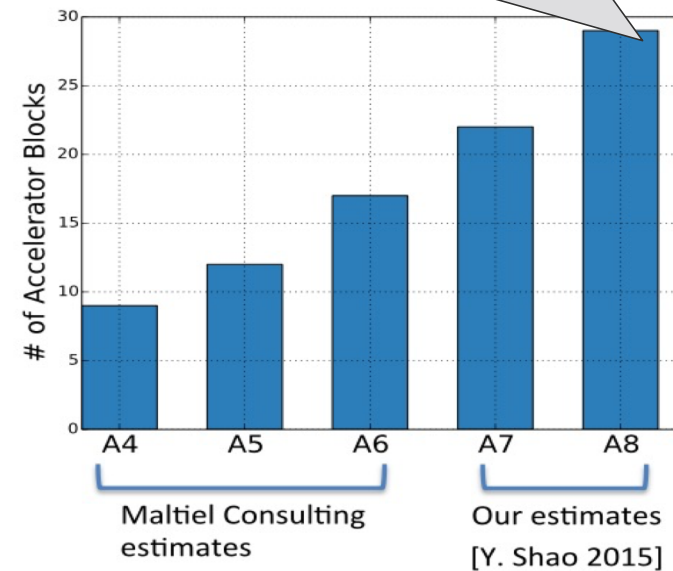


Extreme Specialization Happening Now

(and it will happen to HPC too... will we be ready?)



29 different heterogeneous accelerators in Apple A8 Circa 2016



[www.anandtech.com/show/8562/chipworks-a8]

What is Extreme Heterogeneity?

- Exponentially Increasing parallelism (central challenge for Exascale Computing Project, but will be even worse)
 - Trend: End of exponential clock frequency scaling (end of Dennard scaling)
 - Consequence: Exponentially increasing parallelism
- End of lithography as primary driver for technology improvements
 - Trend: Tapering of lithography scaling
 - Consequence: Many forms of heterogeneous acceleration (not just GPGUs anymore)
- Data movement heterogeneity and increasingly hierarchical machine model
 - Trend: Moving data operands costs more than computation performed on them
 - Consequence: More heterogeneity in data movement performance and energy cost
- Performance heterogeneity
 - Trend: Heterogeneous execution rates from contention and aggressive power management
 - Consequence: Extreme variability and heterogeneity in execution rates

What is Extreme Heterogeneity? (cont'd)

- Diversity of emerging memory and storage technologies
 - Trend: Emerging memory technologies and stall in disk performance improvements
 - Consequence: Disruptive changes to our storage environment
- Increasingly diverse application requirements
 - Trend: Diverse and complex and heterogeneous scientific workflows
 - Consequence: Complex mapping of heterogeneous workflows on heterogeneous systems
- Rapidly expanding community of application developers and users of HPC resources
 - Trend: Larger numbers of domain scientists and non-experts using extreme-scale systems
 - Consequence: Increasing emphasis on productivity and usability

ASCR EH Workshop Charge Letter



Vetter, Jeffrey S.

From: Helland, Barbara <Barbara.Helland@science.doe.gov>
Sent: Thursday, 28 September, 2017 07:58
To: Vetter, Jeffrey S.
Cc: Nowell, Lucy; Susut-Bennett, Cerem; Lee, Steven
Subject: Basic Research Needs for Extreme Heterogeneity Workshop

Jeff,

Thank you for agreeing to be the overall chair for the ASCR workshop focused on the Basic Research Needs for Extreme Heterogeneity. Per our discussions, the workshop will be held on January 23-25, 2018, at the Gaithersburg Marriott Washingtonian Center in Gaithersburg, MD (north of Washington DC). This email confirms ASCR's invitation for you to lead this important ASCR activity.

The workshop will follow the model used by SC's Basic Energy Sciences program with their Basic Research Needs (BRN) workshops. As you know, critical to ASCR's success are the in-person meeting of a broad group of participants from the community and the development of a report that outlines the priority research directions, as identified by the participants, in providing a smart software stack that makes future computers composed of a variety of complex processors and accelerators, new interconnects and deep memory hierarchies productive for future science. The BRN workshops are typically 2.5 days. On the last morning, the panel leads present the priority research directions identified by their panel to the entire group. The afternoon of the third day is reserved for writing by the chairs, panel leads, and other writers who may have been selected by the group.

The charge for the workshop is:

The purpose of this workshop is to identify the priority research directions for ASCR in providing a smart software stack that includes techniques, such as deep learning to make future computers composed of a variety of complex processors, new interconnects and deep memory hierarchies easily used by a broad community of computational scientists. In 2009, the drive to deploy more energy efficient computers led the Oak Ridge Leadership Computing Facility (OLCF) to propose a system upgrade, composed of CPUs coupled with GPUs, which ushered in a new era of heterogeneous high performance computing. The planned Summit upgrade at the OLCF is composed of one CPU coupled to three GPUs and has at least three types of memory. A recent analysis of vendors' current architectural roadmaps is consistent with the increasing heterogeneity that ASCR is seeing in its computing upgrades and indicate that future computers will be more complex and composed of a variety of processing units and accelerators supported by open interconnects and deep memory hierarchies, in other words extremely heterogeneous.

Scientifically, past DOE investments in applied mathematics and computer science basic research and in programs like SciDAC have broadened the community of computational scientists using HPC as one tool to address their grand challenge problems. Nevertheless, significant computer science challenges remain as barriers to efforts to develop a smart software stack that will help increase the usability and programmability of future systems and that will also increase the productivity of the computational scientists. The primary aim for the workshop is to identify the new algorithms and software tools needed from basic research in computer science to enable ASCR's supercomputing facilities to support future scientific and technological advances on SC program's grand challenge problems. ASCR's grand challenges and the resulting priority basic research directions should be identified by spanning existing and next generation computer architectures, including novel technologies that may be developed in the "Post-Moore's Law era" and the promising tools and techniques that are essential to efficient and productive utilization of such architectures. The workshop and subsequent report should define basic research needs and opportunities in computer science research to develop smart and trainable operating and runtime systems, execution models, and

1

programming environments that will make future systems easier to tailor to scientists' computing needs and for facilities to securely deploy.

The chair and co-chairs are responsible for leading the entire workshop planning process. The overall tasks are listed below in approximate chronological order. We will schedule regular conference calls among the chair, co-chairs, and DOE to start the planning process beginning next week.

- Develop the high level workshop structure, including deciding on the number and focus of the panels. Based on the meeting venue, we can have up to 3 panels.
- Based on the panel topics, identify possible plenary topics and speakers.
- Work with DOE to identify panel leads, and then work with the panel leads to identify the workshop participants, including a plan to engage a broad range of DOE Lab personnel, academics and industry representative. Ideally, this plan will provide for inclusion of people who have not participated in ASCR's workshops before. This is a time consuming process that we should begin as soon as possible in order to get the meeting on people's calendars.
- As soon as possible, coordinate preparation of a background document on the status of the field that would be distributed to participants ahead of the workshop. DOE program managers from ASCR will participate in preparing this document.
- During the workshop, synthesize the panels' ideas, guide the identification and definition of priority research directions, and coordinate an oral report to the full workshop at the closing session.
- Critically, coordinate and integrate the topical narratives provided by the panel leads and other identified writers into a final report. As much of the writing as possible is to be completed during the workshop, but follow-up writing is almost always required. ASCR will support a technical editor to help finalize the document.

The goal is to have a final report within 3 months after the workshop in order to maximize the report's impact on programmatic planning.

We really appreciate your willingness to lead this essential planning activity for ASCR.

Barbara Helland
Associate Director
for Advanced Scientific Computing Research
Office of Science

Barbara Helland
Associate Director, Advanced Scientific Computing Research
Office of Science
Department of Energy
Phone: 301-903-7486
Email: Barbara.Helland@science.doe.gov



ASCR Extreme Heterogeneity Workshop

- January 23-25, 2018
- Virtual workshop - face-to-face meeting canceled due to government shutdown
- Several plenary talks on hardware trends, memory technology, quantum computing, machine learning, workflow
- Attendees chosen based on submitted white papers
- Breakout groups
 - Programming Environments, Models, and Languages
 - Data Management and I/O
 - Data Analytics and Workflows
 - Operating Systems and Resource Management
 - Software Development Methodologies
 - Modeling and Simulation for Hardware Characterization
 - Programming Environments: Compilers, Libraries, and Runtimes
 - System Management, Administration, and Job Scheduling
 - Crosscut: Productivity, Composability, Interoperability
 - Crosscut: Portability, Code Reuse, and Performance Portability
 - Programming Environments: Debugging, Autotuning, Specialization
 - Crosscut: Resilience and Power Management
- <https://www.ornl.gov/ExHeterogeneity2018>

EH Workshop Organizing Committee

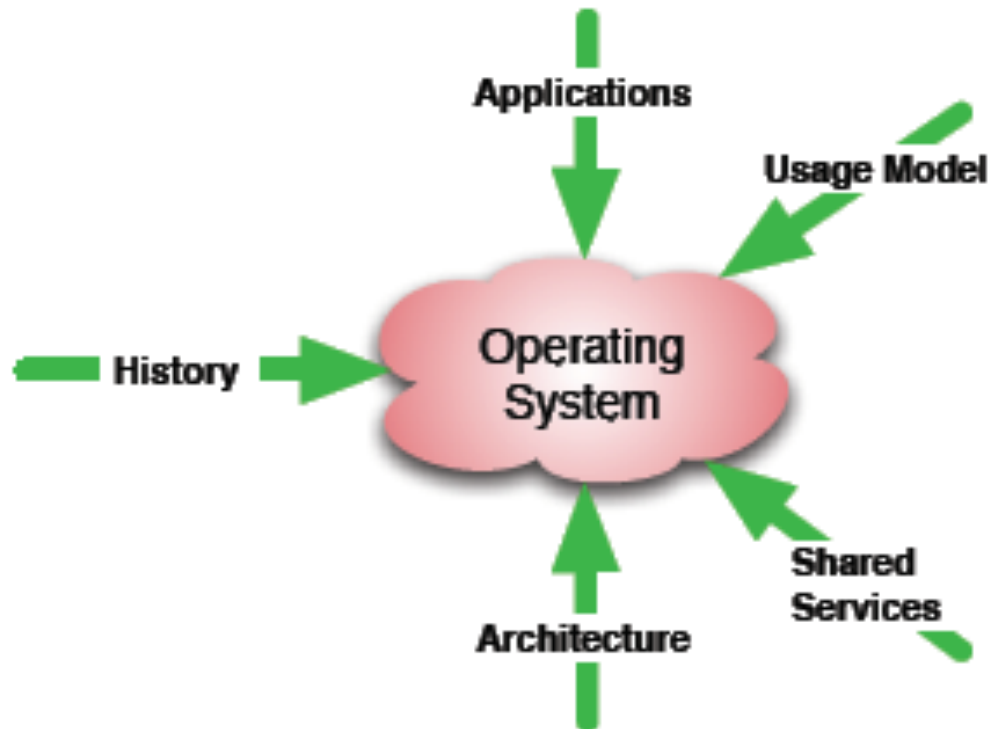
- Jeffrey Vetter, Chair (ORNL)
- Katie Antypas (LBNL)
- Ron Brightwell (SNL)
- David Donofrio (LBNL)
- Maya Gokhale (LLNL)
- Travis Humble (ORNL)
- Pat McCormick (LANL)
- Rob Ross (ANL)
- Catherine Schuman (ORNL)
- John Shalf (LBNL)
- Brian Van Essen (LLNL)
- Shinjae Yoo (BNL)

Program Manger: Lucy Nowell

Safe Harbor Statement

- This is my view on the research directions and priorities resulting from the workshop
- The final report is currently in development and my views may or may not be reflected in the report

Factors Influencing OS Design



Architecture

- System-on-Chip (SoC)
 - Hardware specialization
 - OS/R needs to be aware of custom hardware capabilities
 - Potentially large collection of hardware capabilities where only a few may be used at a time
 - A single node will not be a single cache-coherent physical address space (true today)
- Photonic interconnects
 - Load/store across a larger domain
 - More intelligent memory controllers
 - Perhaps programmable by OS/R or application
 - Converged with network interface
 - Nodes will look more like racks, racks will look more like systems
- Special-purpose systems will become more general
 - OS will have to be engineered to adapt more easily to new hardware
- Trust model will have to evolve
 - Security model for users and applications likely needs to change
- OS will become much more distributed

Applications

- Increased complexity
 - Reduce complexity through abstractions, componentization, and composition
 - Decompose applications into tasks and services
 - OS/R will need to provide mechanisms for service discovery and composition
- Access to system services
 - Traps and blocking system calls are already insufficient
 - Convergence between OS and RTS
 - Expose hardware directly to application
- Tools are applications too
 - Tools typically depend more on system services
 - Less human interaction with tools
 - Consumer of diagnostic and debugging information may be the OS or RTS
- Rethink the connections between OS/R and programming environment
- Likely to be event-driven at some level

Usage Model

- Need to move beyond batch-scheduled, space-shared, non-interactive jobs
 - Dedicated resources versus shared resources
 - More interactivity with users and application services
 - Need to develop a new cost charging model for facilities
- Implicit versus explicit allocation and management of resources
 - Already seeing limitations with explicitly allocating cores, nodes, memory (burst buffers) etc.
 - OS/R will likely need to determine resources implicitly and be elastic
 - Methods for handling resource failures
- Data-centric versus compute-centric view of system
 - Differentiating between HPC and Cloud/BigData approaches
- Support new methods of moving data on and off of the system

Shared Services

- RAS System (Reliability/Availability/Serviceability)
 - Instrumentation and analysis
 - System health monitoring
 - In-band and/or out-or-band
 - Global Information Bus
- External resources
 - External connectivity to network and storage
 - Streaming data from external instruments
 - New methods of data ingest/egest

History

- Legacy programming interfaces
 - POSIX probably needs to go away for more than just I/O
 - Glibc may not be the RTS of the future
 - Requires support for incremental adoption
- Standard protocols
 - Identify abstraction layers that allow for evolution
- May finally have to move away from Unix model
 - Convergence of memory and storage is a fundamental change for the OS
 - Everything is really not a file
- Need to balance between starting from scratch and supporting existing infrastructure

Pop Quiz

- Which of the following is the most important factor in determining whether or not a technology is adopted by Sandia application developers?
 - Performance
 - Scalability
 - Maturity
 - Sustainability
 - Ease of integration
 - Testability

Sustainability

- Application developers will more readily adopt a technology that they know will be sustained over time
- This is especially true with new C++ language features

EH Priority Research Directions

- Maintaining and improving programmer productivity
 - Flexible, expressive, programming models and languages
 - Intelligent, domain-aware compilers and tools
 - Composition of disparate software components
- Managing resources intelligently
 - Automated methods using introspection and machine learning
 - Optimize for performance, energy efficiency, and availability
- Modeling & predicting performance
 - Evaluate impact of potential system designs and application mappings
 - Model-automated optimization of applications
- Enabling reproducible science despite non-determinism & asynchrony
 - Methods for validation on non-deterministic architectures
 - Detection and mitigation of pervasive faults and errors
- Facilitating Data Management, Analytics, and Workflows
 - Mapping of science workflows to heterogeneous hardware and software services
 - Adapting workflows and services to meet facility-level objectives through learning approaches

Reduce time to verifiable discovery despite diverse application domains and an exponential increase in architectural complexity from rapidly changing heterogeneous systems:

Managing Resources Intelligently (1/3)

- OS/RTS Design: Hardware resources will become more complex and diverse. The operating system (OS) and runtime system (RTS) must integrate special-purpose devices and accelerators. The OS cannot assume all resources on a node are identical and dedicated devices
 - OS/RTS must be efficient and sustainable for an increasingly diverse set of hardware components
 - Must provide capability for dynamic discovery of resources as power/energy constraints impose restrictions on availability

Managing Resources Intelligently (2/3)

- Decentralized resource management: New scalable methods of coordinating resources must be developed that allow policy decisions and mechanisms to co-exist throughout the system. Hardware resources are becoming inherently adaptive, making it increasingly complex to understand and evaluate optimal execution and utilization
 - System software must be enhanced to coordinate resources across multiple levels and disparate devices in the system
 - Upper layers can't assume ownership of all resources
 - Must leverage cohesive integration of performance introspection and programming system abstractions to provide more adaptive execution
 - Optimization without human intervention
 - Improved information flow between application and resource management system

Managing Resources Intelligently (3/3)

- Autonomous resource optimization: Responsibility for efficient use of resources must shift from the user to the system software; must employ sophisticated and intelligent approaches optimize selection of resources to application needs
 - Need more automated methods using machine learning to optimize the performance, energy efficiency, and availability of resources for integrated application workflows
 - Implicit rather than explicit allocation and management of resources
 - More sophisticated usage models beyond batch-scheduled, spaced-shared nodes adds significant complexity to the management of system resources
 - Map the machine to the application rather than vice-versa

Exploring Potential Solutions

- “We’ll explore using that feature when it gets into the Standard.”
- “We’ll consider task-based models when it’s clear MPI everywhere won’t work anymore.”
- That’s too late
- Applications must continuously engage in exploration and evaluation of technologies
- Otherwise potential solutions won’t be available when they are required
- Must provide methods of incremental adoption and evaluation

Issues/Concerns

- Managing the memory hierarchy
 - Lots of evidence that the RTS/OS are not good at this for HPC
- Increasing complexity and responsibility of the RTS
 - Pushing complexity to the RTS with less info
- Resource requirements of the RTS
 - Potentially significant overhead
- Compelling application evaluation
 - Applications need to exercise the advanced RTS functionality
 - Implementation bias
- Application performance portability
 - From laptop to beyond exascale
- Transparency is in the eye of the application developer
 - Need to support both experts and ambivalent

More Issues/Concerns

- Cost of modularity
 - Not all RTS services should be componentized
- Ability to constrain the problem
 - Too many hardware and application “knobs”
- Performance portability of the RTS
 - Not any easier to solve than application performance portability
- Dependence on hardware advancements
 - Inability to demonstrate compelling results on current systems
- Lack of standard low-level network API
 - Fundamental issue for RTS communication
- HPC market pressure
 - Influence of non-HPC “solutions”

Final Issues/Concerns

- Amount of asynchrony
 - Ability of algorithms to reduce global operations
- Mechanisms to support event-driven capability
 - More efficient ways to enable adaptivity
- Walking before running
 - Make progress at small scale while working towards large scale
- Programming system evaluation and comparison
 - Need scientific approach to measuring effectiveness of different programming systems

Acknowledgments

- EH Workshop Organizing Committee
- Participants in the OS/Resource Management breakout sessions
- Slides 6-10 were stolen from John Shalf

Questions or comments?