# Hobbes Extreme Scale OS

## Ron Brightwell

R&D Manager

Scalable System Software

Center for Computing Research

Sandia National Laboratories

SPPEXA Workshop on Application Interfaces for an Exascale OS

TU Dresden

December 8, 2014

FFMK Workshop

**Sandia National Laboratories**

*Exceptional service in the national interest*
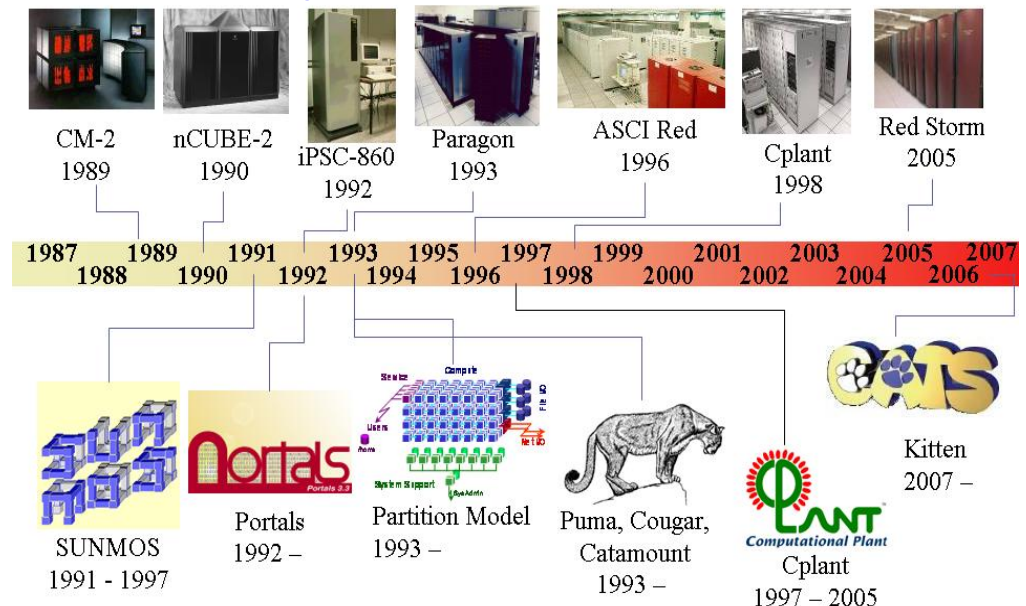
# System Software@Sandia

- Established the functional partition model for HPC systems
  - Tailor system software to function (compute, I/O, user services, etc.)
- Pioneered the research, development, and use of lightweight kernel operating systems for HPC
  - Only DOE lab to deploy OS-level software on large-scale production machines
  - Provided blueprint for IBM BG/L,Q CNK
- Set the standard for scalable parallel runtime systems for HPC
  - Fast application launch on tens of thousands of processors
- Significant impact in the design and of scalable HPC interconnect APIs
  - Only DOE lab to deploy low-level interconnect API on large-scale production machines
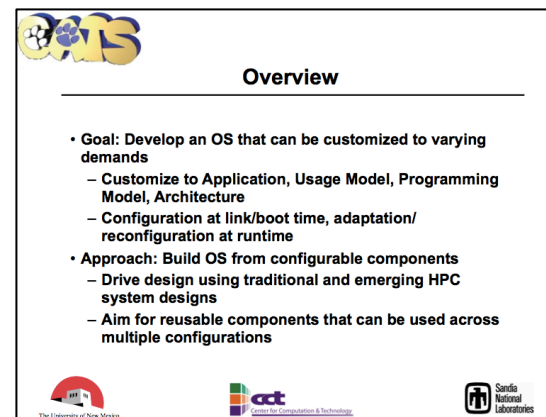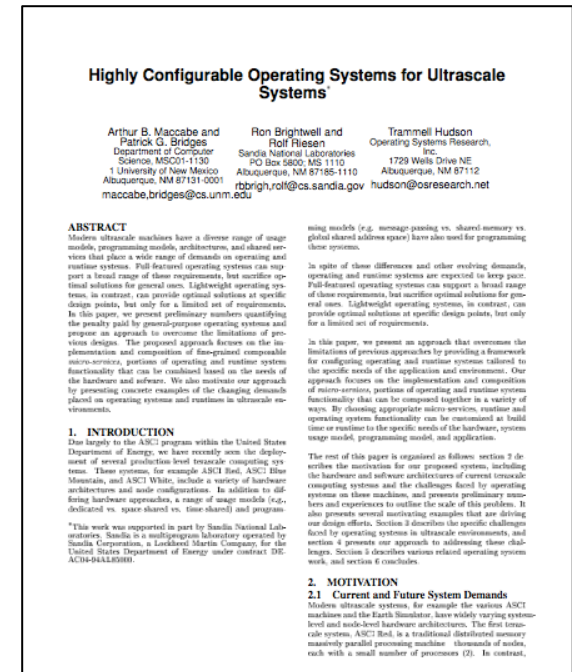
**AWARDS:**
- **1998** Sandia Meritorious Achievement Award, TeraFLOP Computer Installation Team
- **2006** Sandia Meritorious Achievement Award, Red Storm Design, Development and Deployment Team
- **2006** NOVA Award Red Storm Design and Development Team
- **2009** R&D 100 Award for Catamount N-Way Lightweight Kernel
- **2010** Excellence in Technology Transfer Award, Federal Laboratory Consortium for Technology Transfer
- **2010** National Nuclear Security Administration Defense Programs Award of Excellence

# Configurable OS Project

- **Part of the US DOE FAST-OS Program**
  - Partnership between Sandia, UNM, LSU
  - 2005-7
- **Lessons learned**
  - HPC users want Linux environment
    - Most only care about toolchain
  - Very little middle ground between LWK and Linux
  - Difficult to concentrate on $500K research project when $75M Red Storm system needs attention

# Virtualization May Help

# Kitten Lightweight Kernel

- Monolithic, C code, GNU toolchain, Kbuild configuration

- Supports x86-64 architecture only, porting to ARM
  - Boots on standard PC architecture, Cray XT, and in virtual machines
  - Boots identically to Linux (Kitten bzImage and init_task)

- Repurposes basic functionality from Linux
  - Hardware bootstrap
  - Basic OS kernel primitives (lists, locks, wait queues, etc.)
  - Directory structure similar to Linux, arch dependent/independent dirs

- Custom address space management and task management
  - User-level API for managing physical memory, building virtual address spaces
  - User-level API for creating tasks, which run in virtual address spaces

- Small, highly reliable code base

- Focused on scalable HPC applications
  - Low noise
  - Small memory footprint

- Open source and freely available

# Palacios Virtual Machine Monitor

- OS-independent embeddable virtual machine monitor
  - Can be combined with Kitten or Linux

- Full system virtualization
  - No need to modify guest OS

- Supports running multiple guests concurrently

- Makes extensive use of virtualization extensions in modern Intel and AMD x86 processors

- Passthrough resource partitioning

- Extensive configurability

- Low noise

- Open source and freely available

- Small, highly reliable code base

- Developed at Northwestern and University of New Mexico

# Low Overhead for Palacios+Kitten on Red Storm (2009)



CTH: multi-material, large deformation, strong shockwave simulation

# DOE/ASCR X-Stack 2012

# Enabling Exascale Hardware and Software Design through Scalable System Virtualization

Patrick G. Bridges, University of New Mexico; Peter Dinda, Northwestern University; Jack Lange, University of Pittsburgh; Kevin Pedretti, Sandia National Laboratories; Stephen Scott, Oak Ridge National Laboratory

## Overview

In this project, we are investigating system software tools to accelerate the development and use of exascale systems. In particular, we are developing new system virtualization techniques to enable the development of the hardware and software innovations needed to enable exascale systems. Virtualization techniques provide traction on a wide range of exascale design and development issues, as described below.

We are using virtualization to enable the design, development, and use of exascale systems. Virtualization allows new hardware and software features to be prototyped as extensions to a virtual machine monitor (VMM), making them immediately available for experimentation and use. Furthermore, it allows new system software and runtime stacks to be launched above production host operating systems without the need for dedicated system time. The VMM

# DOE LAB 13-02 FOA – January 2013

## Exascale Operating and Runtime Systems Program

- $6M of funding for OS/R research at US DOE labs

- Focus areas

  - Power management
    - Adaptive power management to meet 20 MW goal

  - Support for dynamic programming environments
    - Manage billions of threads

  - Programmability and tuning support
    - Dynamic adaptation and debugging

  - Resilience
    - Predict, detect, contain, and recover from faults

  - Heterogeneity
    - Hierarchical process and memory systems

  - Memory management
    - Use of new memory technologies

  - Global optimization
    - Manage resources with a system-wide view

# Exascale OS/R Focus is on Hardware

- Reliability/Resilience
- Power/Energy
- Heterogeneity
- Memory hierarchy
- Cores, cores, and more cores

- Risk
  - Hardware advancements and investments can provide orders of magnitude improvement
  - OS/R advancements can provide double-digit percentage improvement

# What About Applications?

- Focus is on parallel (many-core) programming model
  - Adaptive runtime systems
  - Node-level resource allocation and management
  - Managing locality
  - Extracting parallelism
  - Introspective, adaptive capabilities
    - US DOE/ASCR XPRESS project is addressing OS support for adaptive RTS

- Risk
  - Incremental approach (OpenMP) wins
    - Advanced runtime capabilities are overkill
  - No clear on-node parallel programming model winner
    - Difficult to optimize OS/R

# Application Composition Will Be Increasingly Important at Extreme-Scale

- More complex workflows are driving need for advanced OS services and capability
  - Exascale applications will continue to evolve beyond a space-shared batch scheduled approach
- HPC application developers are employing ad-hoc solutions
  - Interfaces and tools like mmap, ptrace, python for coupling codes and sharing data
- Tools stress OS functionality because of these legacy APIs and services
- More attention needed on how multiple applications are composed
- Several use cases
  - Ensemble calculations for uncertainty quantification
  - Multi-{material, physics, scale} simulations
  - In-situ analysis
  - Graph analytics
  - Performance and correctness tools
- Requirements are driven by applications
  - Not necessarily by parallel programming model
  - Somewhat insulated from hardware advancements

# Hobbes Project Goals

- Deliver prototype OS/R environment for R&D in extreme-scale scientific computing

- Focus on application composition as a fundamental driver
  - Develop necessary OS/R interfaces and system services required to support resource isolation and sharing
  - Support complex simulation and analysis workflows

- Provide a lightweight OS/R environment with flexibility to build custom runtimes
  - Compose applications from a collection of enclaves

- Leverage Kitten lightweight kernel and Palacios lightweight virtual machine monitor
  - Enable high-risk high-impact research in virtualization, energy/power, scheduling, and resilience

- Enable High-Risk/High-Impact R&D in key areas

# Fundamental Principles of Our Approach

- OS/R must be viewed as technologies that enable and support the research and development of other critical capabilities required for effective use of extreme-scale high-performance computing (HPC) systems

- OS/R support for composition of applications is a critical capability that will be the foundation of the way extreme-scale systems must be used in the future

- Addressing near-term OS/R challenges such as energy/power, scheduling, and resilience without considering application composition will lead to incomplete solutions

# About the Name….

# Or Possibly…

- HPC
- OS
- Building
- Blocks for
- Extreme-scale
- Systems

# A Deeper Look at Composition

## Intra-Node Composition

- Components co-located on same set of nodes
- Virtualization used to isolate NOS environments on each node
- Composition (coupling) takes place via shared memory

## Inter-Node Composition

- Components deployed to separate sets of nodes
- Composition (coupling) takes place via network



Enclave OS/Rs

Node OS/Rs

Node Virtualization Layers

Physical Nodes

# Composition of Enclaves

- An enclave provides a single OS/R environment to the application
- Hobbes approach is to provide the minimum "amount" of OS/R required by the application (do what is necessary and get out of the way!)
- Modern, complex applications are increasingly created by assembling (often substantial) software components
  - E.g., analytics connected to applications, code coupling, application frameworks, …
- Components may have distinct requirements for OS/R support
- Two options:
- Assemble an all-in-one OS/R stack that satisfies all component needs
  - Potential challenges at both OS and RTS levels
  - Requires integration work for every combination supported
- Provide each component the OS/R it needs, and provide efficient, low-level mechanisms to connect the components (and the OS/Rs)

# Hobbes Node Architecture
## Independent Operating and Runtime Systems

Policies to manage the VMs on a single node.

**Global Information Bus**

**On Node Management**

EOS 1

EOS 2

SGOS

VM 1

Application 1

RT 1

NOS 1

VM 2

Application 2

RT 2

NOS 2

VMs can share the resources via time sharing or space sharing. This is managed by the SGOS

User

Kernel

**VM Management Module**

**HAL (Hardware Virtualization)**

Additional mechanisms needed to manage multiple VMs. Run in kernel mode to take advantage of VM support in modern processors (Palacios).

Basic mechanisms needed to virtualize hardware resources like address spaces (Kitten).

# Hobbes Node Architecture
## Unified Operating and Runtime Systems

Policies to manage the VMs on a single node.

Global Information Bus

On Node Management

EOS

SGOS

VM

Application 1

Application 2

Unified RT

NOS

Sharing among applications is managed by the NOS and Unified RT. The SGOS has a minimal role.

User

Kernel

VM Management Module

HAL (Hardware Virtualization)

Additional mechanisms needed to manage multiple VMs. Run in kernel mode to take advantage of VM support in modern processors (Palacios).

Basic mechanisms needed to virtualize hardware resources like address spaces (Kitten).

# Hobbes Node Architecture
## HPC Application and Runtime

Policies to manage the VMs on a single node.

Global Information Bus

On Node Management

EOS

SGOS

VM

Application

LW RT

No sharing of node resources. The on-node GOS is minimal and the VM module might be gone.

User

Kernel

VM Management Module

HAL (Hardware Virtualization)

Additional mechanisms needed to manage multiple VMs. Run in kernel mode to take advantage of VM support in modern processors (Palacios).

Basic mechanisms needed to virtualize hardware resources like address spaces (Kitten).

# Composition Examples

- SNAP + Analytics
  - "SNAP calculates synonymous and non-synonymous substitution rates based on a set of codon-aligned nucleotide sequences." (HIV related)
  - Proxy app from LANL used for example
- GTC-P + Analytics
  - Fusion simulation testing/proxy app used to test new hardware and algorithm integration into the PIC model. (PPPL)
  - Analytics generate statistics on particles (histograms), sorts, and filters on bounding boxes
- LAMMPS + Analytics
  - Full, production molecular dynamics application from Sandia
  - Analytics look for crack formation by calculating atomic spacing in output data to change simulation from coarse to fine grained.

# Hobbes Has Several Components

- <span style="color:green">Node Virtualization Layer</span>
- <span style="color:green">Enclave OS</span>
- Scheduling
- Programming Models
- Global Information Bus
- Resilience
- Power/Energy

# Hobbes Team

| Institution | Person | Role |
|---|---|---|
| Georgia Institute of Technology | Karsten Schwan | PI |
| Indiana University | Thomas Sterling | PI |
| Los Alamos National Lab | Mike Lang | PI |
| Lawrence Berkeley National Lab | Costin Iancu | PI |
| North Carolina State University | Frank Mueller | PI |
| Northwestern University | Peter Dinda | PI |
| Oak Ridge National Laboratory | David Bernholdt | PI |
| Oak Ridge National Laboratory | Arthur B. Maccabe | Chief Scientist |
| Sandia National Laboratories | Ron Brightwell | Coordinating PI |
| University of Arizona | David Lowenthal | PI |
| University of California – Berkeley | Eric Brewer | PI |
| University of New Mexico | Patrick Bridges | PI |
| University of Pittsburgh | Jack Lange | PI |

# Node Virtualization Layer (NVL)

# Why a Node Virtualization Layer?

- Flexibility, support multiple OS/R stacks simultaneously
  - There is likely to be no one-size-fits-all OS/R stack, lots of exploration
  - Co-location of VMs, efficient sharing of resources between enclaves
  - Native environment freed from legacy constraints
- Low overhead
  - Our past work has shown CPU and memory overheads negligible
  - Network I/O is still an issue, but tractable
- Industry momentum
  - Virtualization has been commoditized, is everywhere
  - Academic and student mindshare, where the jobs are
- Mostly orthogonal to "FusedOS" approach others are taking
  - FusedOS could run in NVL VM or natively, in the same machine
  - NVL could be co-designed with FusedOS
- Already doing node OS R&D under XPRESS

# Exploring Spectrum of Virtualization

- Virtualization doesn't have to be "big and heavy"
  - Don't have to trap everything
  - VMM can setup paths to hardware, then get out of way
- There are multiple virtualization architectures, not just one
  - Hobbes NVL team working across spectrum (Blue items, research in Light Blue)

**Lightest Weight**

**LWK**
Virtual Linux
Evironment
(**Kitten**, CNK)

**Combo OS**
Multiple-native
OSes
(**Pisces**, FusedOS)

**Full HW VM**
Runs Unmodified
Guest OSes, Passthru
(**Palacios**, KVM, …)

**Heaviest Weight**

**LWK**
Custom
(Catamount,
**HybridVM**)

**Para-virtual**
Implicit,
VMM Changes
Guest OS
(**Gears, Guarded
Modules**)

**Para-virtual**
Explicit,
Guest OS Modified
or Augmented
(Orig. Xen,
**Device Drivers**)

**Software Virt**
Emulate HW, Binary
Translation, …
(Qemu, Vmware,
**Emulate HW Trans
Memory pre-product**)

# NVL Has Multiple Levels of Virtualization

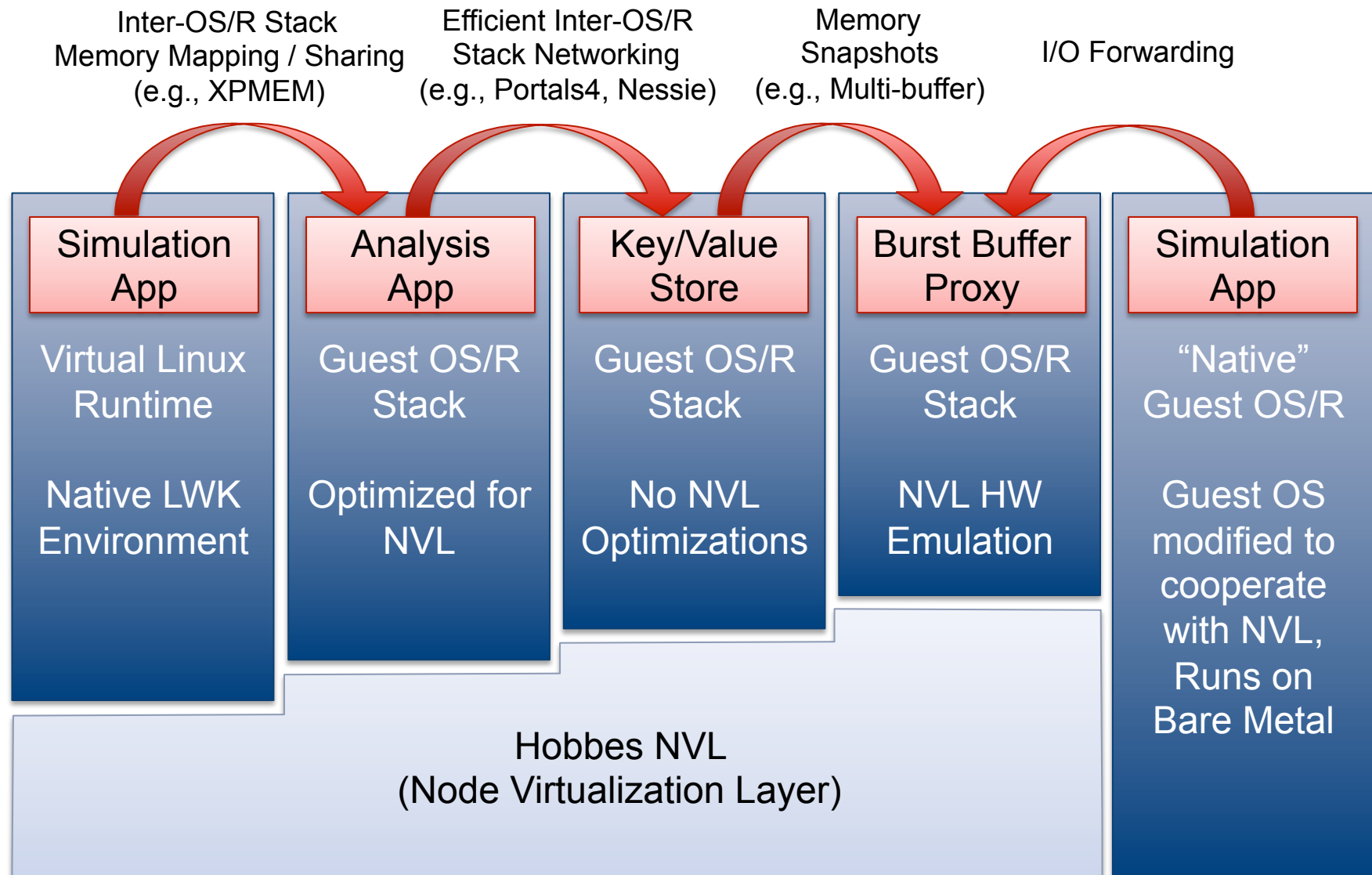- Existing Hypervisors typically support one level
- NVL couples LWK "native" runtime with guest OS/R stacks

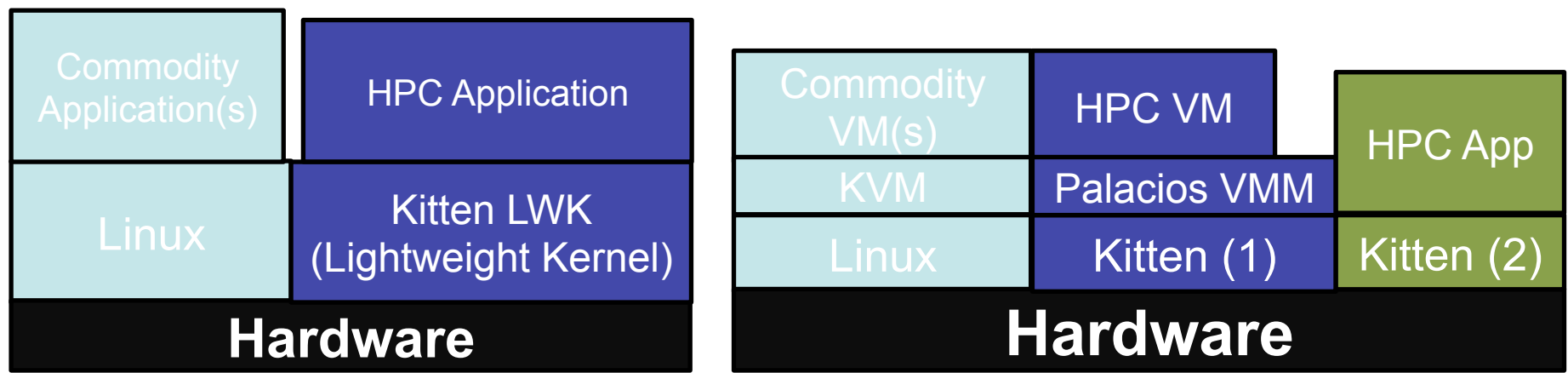| Application | Application | Application | Application | Application |
|---|---|---|---|---|
| Virtual Linux Runtime | Guest OS/R Stack | Guest OS/R Stack | Guest OS/R Stack | "Native" Guest OS/R |
| Native LWK Environment | Optimized for NVL | No NVL Optimizations | NVL HW Emulation | Guest OS modified to cooperate with NVL, Runs on Bare Metal |

Hobbes NVL
(Node Virtualization Layer)

# NVL Provides Composition Mechanisms

Inter-OS/R Stack
Memory Mapping / Sharing
(e.g., XPMEM)

Efficient Inter-OS/R
Stack Networking
(e.g., Portals4, Nessie)

Memory
Snapshots
(e.g., Multi-buffer)

I/O Forwarding

| Simulation App | Analysis App | Key/Value Store | Burst Buffer Proxy | Simulation App |
|---|---|---|---|---|
| Virtual Linux Runtime<br><br>Native LWK Environment | Guest OS/R Stack<br><br>Optimized for NVL | Guest OS/R Stack<br><br>No NVL Optimizations | Guest OS/R Stack<br><br>NVL HW Emulation | "Native" Guest OS/R<br><br>Guest OS modified to cooperate with NVL, Runs on Bare Metal |

Hobbes NVL
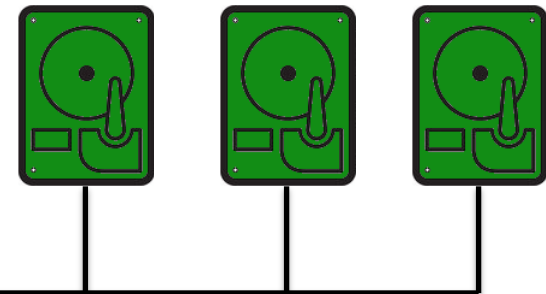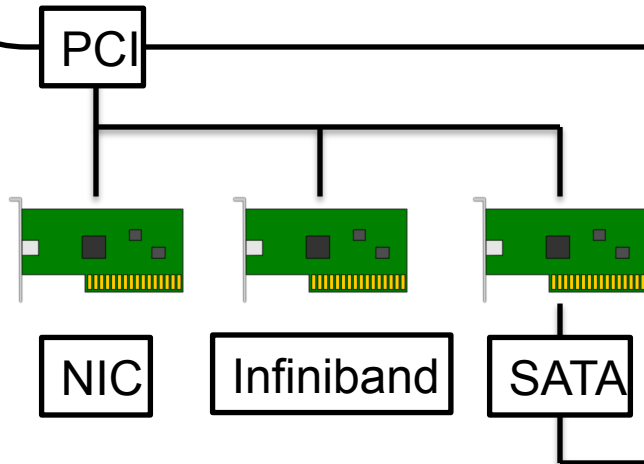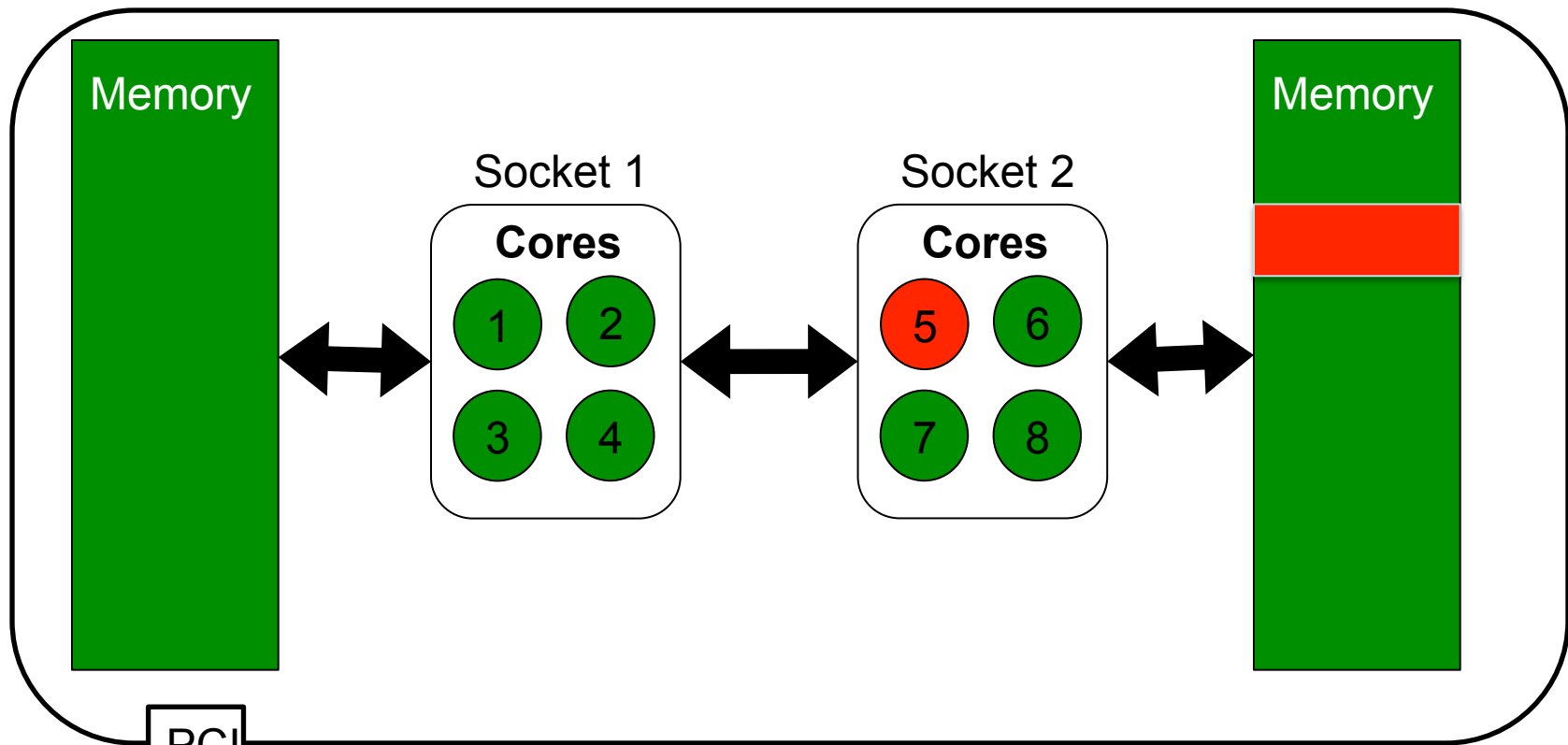(Node Virtualization Layer)

# Pisces Multi-stack Architecture is Our Starting Point

- Goals
  - Fully isolated and independent operation
  - OS Bypass communication
  - No cross-kernel dependencies
  - Leverage cloud platforms
- Uses Linux hot un/plug to bring cores, memory, and devices offline
- Recent modifications to Kitten
  - Boot process that initializes only subset of offline resources
  - Dynamic resource (re)assignment to Kitten
  - Cross stack shared memory communication (XPMEM)
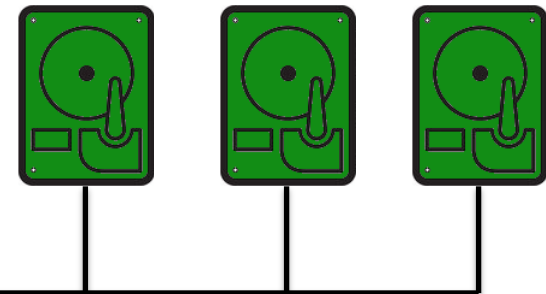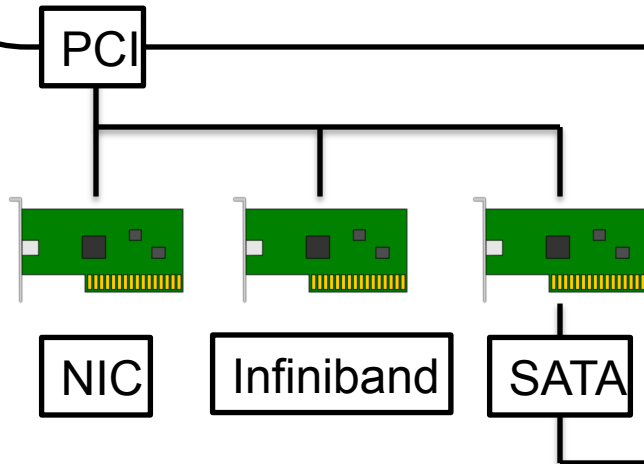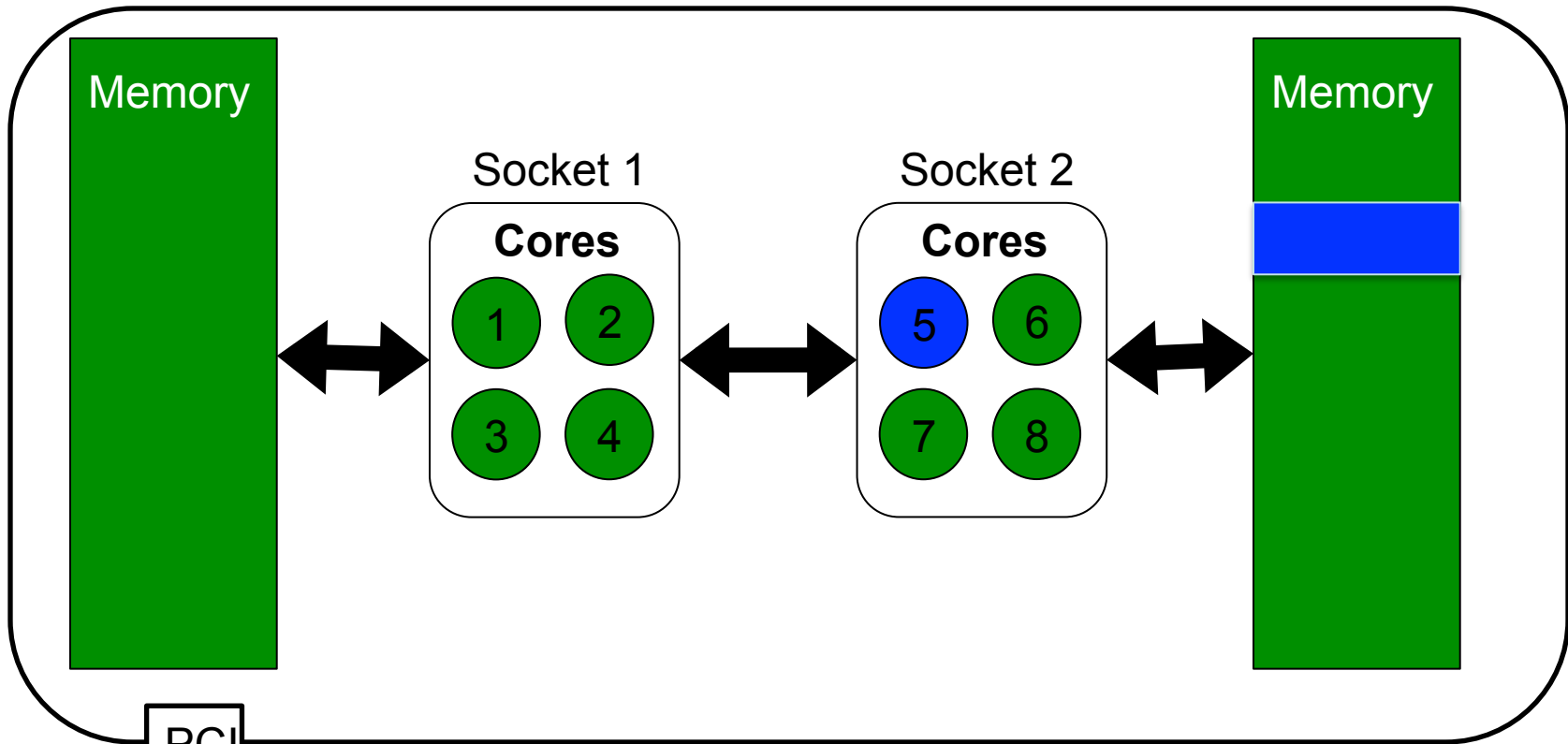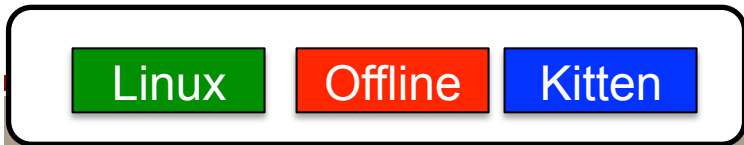  - Block Driver Interface

Memory

Socket 1

**Cores**

1  2

3  4

Socket 2

**Cores**

5  6

7  8

Memory

PCI

NIC   Infiniband   SATA

Linux   Offline   Kitten

Memory

Memory

Socket 1

Socket 2

**Cores**

**Cores**

1 2

5 6

3 4

7 8

PCI

NIC

Infiniband

SATA

Linux Offline Kitten

Memory

Socket 1

Socket 2

Memory

**Cores**

1 2
3 4

**Cores**

5 6
7 8

PCI

NIC

Infiniband

SATA

Linux   Offline   Kitten

Memory

Socket 1
**Cores**
1  2
3  4

Socket 2
**Cores**
5  6
7  8

Memory

PCI

NIC  Infiniband  SATA

Linux  Offline  Kitten

Memory

Socket 1

Cores

1 2
3 4

Socket 2

Cores

5 6
7 8

Memory

PCI

NIC  Infiniband  SATA

Linux  Offline  Kitten

Memory

Socket 1

**Cores**

1 2
3 4

Socket 2

**Cores**

5 6
7 8

Memory

PCI

NIC | Infiniband | SATA

Linux | Offline | Kitten

FFMK Workshop

Memory

Memory

Socket 1

Socket 2

**Cores**

**Cores**

1 2
3 4

5 6
7 8

PCI

NIC

Infiniband

SATA

Linux    Offline    Kitten

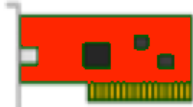Memory

Memory

Socket 1

Socket 2

**Cores**

1 2

3 4

**Cores**

5 6

7 8

PCI

NIC

Infiniband

SATA

Linux    Offline    Kitten

# Node Virtualization Layer Status

- Pisces multi-stack architecture tools implemented and functional

  

  Co-Kernel Architecture

  - Host boots Linux
  - Cores and memory can be taken from Linux, forming one or more containers
  - Kitten can be launched in each container
  - Each Kitten instance operates cooperatively with Linux as a co-kernel
  - Each co-kernel can run a different application
    - Or guest OS via Palacios
  - Containers can be dynamically resized without rebooting
    - Number of cores and size of memory can grow and shrink
- Ported to Cray Linux Environment
- Multi-enclave launch working on Cray XK7 testbed at Sandia

# NVL Status (continued)

- **Using XPMEM for inter-OS shared memory**
  - XPMEM allows address space sharing between distinct processes
  - The OS running on an NVL instance can export and attach to memory regions from co-kernels running on the same NVL
  - All combinations working: Linux <->Linux, Linux<->Kitten, Kitten<->Kitten (where Linux is either Host Linux or Guest Linux)

- **Transparently Consistent Asynchronous Shared Memory (TCASM)**
  - Allows for asynchronously exporting a snapshot of a memory region to many observers
  - Completed initial port from Linux to Kitten

# Pisces Single-Node Performance Evaluation

- Dell R450 server
    - Two six-core Intel IvyBridge Xeons (12 cores total)
    - 24 GB RMA in two NUMA domains
- CentOS 7 Linux Distribution
    - Linux kernel v3.16
- Benchmarks
    - miniFE from Mantevo mini-app suite
    - HPC Challenge RandomAccess Benchmark
- OS environments
    - CentOS – CentOS Linux on all 12 cores
    - Kitten/KVM – CentOS Linux on 6 cores, Kitten KVM guest on 6 cores
    - Co-Kernel – CentOS Linux on 6 cores, Kitten co-kernel on 6 cores
- Execution environment
    - Without background noise
        - Benchmark executed alone in one NUMA domain
    - With background noise
        - Benchmark executed in one NUMA domain, kernel compile in one NUMA domain

# Co-Kernel Has Better Performance and Performance Isolation

**miniFE**

**RandomAccess**

# TCASM Performance



Kitten TCASM implementation showing good performance and scalability
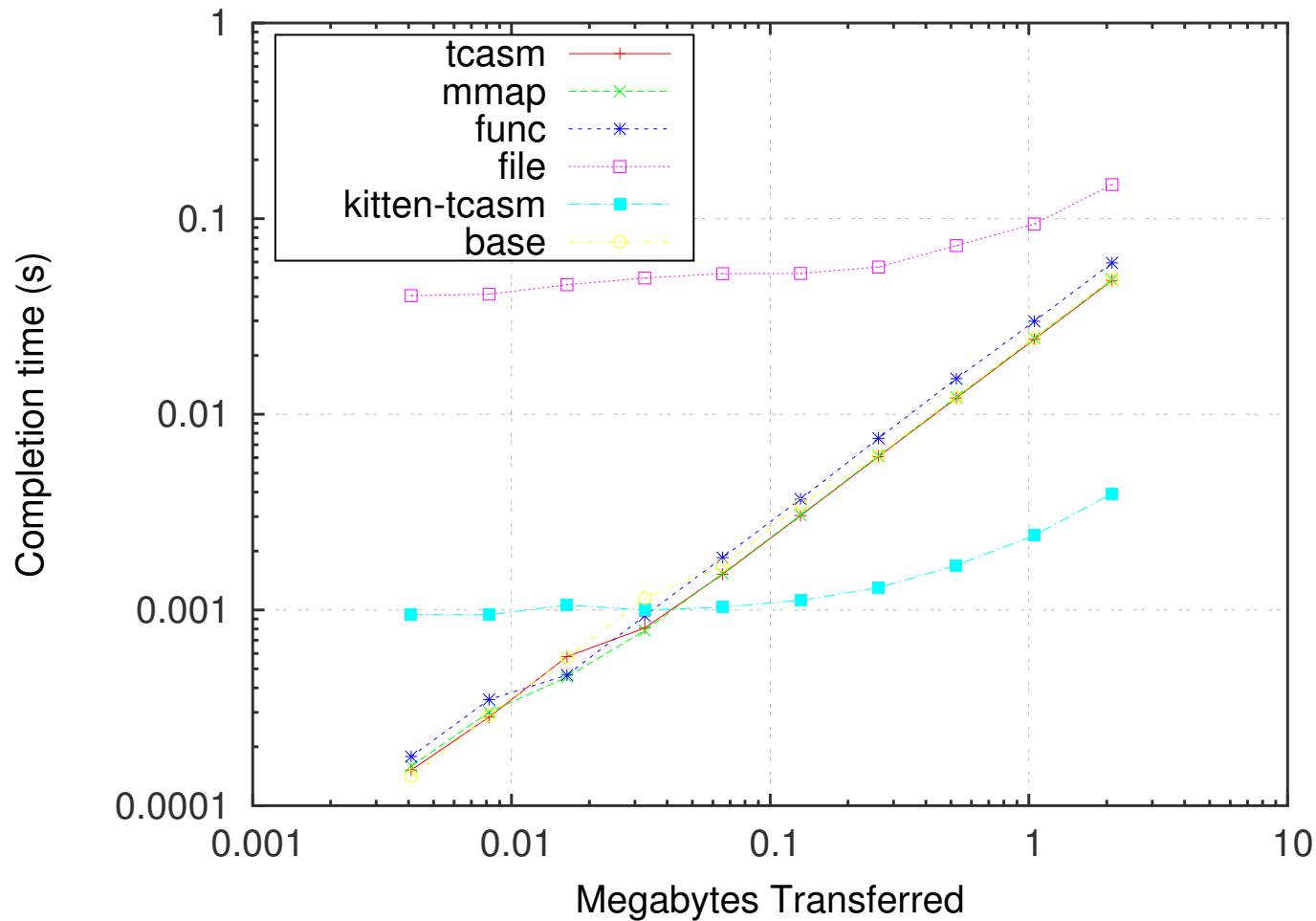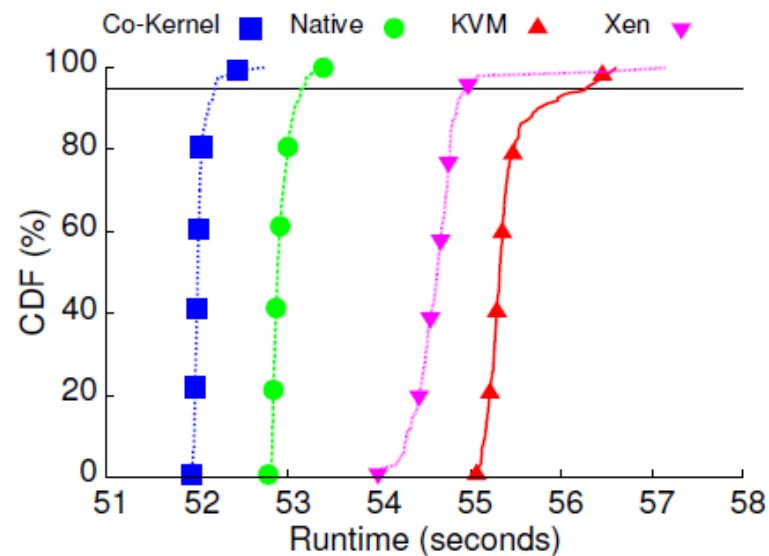
Results are for producer/consumer microbenchmark, both processes running in Kitten

SNAP running in Kitten coupled with analytics running in Linux is under test

# Pisces Multi-Node Performance Evaluation

- 8 node InfiniBand cluster
- Each node has 16 cores, 24 GB of memory, and two HDDs
- Hot unplugged 8 cores, 12 GB (one socket), and one HDD from Linux on each node
- Used Pisces to launch Kitten/Palacios on 8 nodes
  - Kitten SATA driver manages one HDD on each node for local disk I/O for Palacios VMs
- Each node runs
  - A cloud-serving Ubuntu guest on a KVM host
    - Runs Hadoop data nodes with a machine learning benchmark (Mahout)
    - Bridged to a 1 GigE NIC connect to GigE switch
  - An HPC-serving Fedora 19 guest running on either same KVM host or on Palacios/Kitten
    - HPCCG and Cloverleaf mini-apps
    - Connected to passthrough IB network
- Measured the cumulative distribution function on a few hours (300+ runs) of HPCCG running on both KVM and Palacios/Kitten while simultaneously running Hadoop on the same node

# Co-Kernel Has Better Performance and Performance Isolation



(a) HPCCG

(b) CloverLeaf

# NVL Plan

- ADIOS working over XPMEM
  - ADIOS provides higher-level interface for coupling than using XPMEM directly
  - Application composition use cases initially targeting ADIOS

- Develop NVL name service
  - Will allow NVL OS instances running on the same node to post and query available resources that can be shared for composition
  - Eventually may become part of the Global Information Bus

- TCASM
  - Need to better understand initial performance evaluation
  - Implement inter-OS TCASM, integrate with XPMEM
  - TCASM will be a new type of XPMEM memory region
  - Will compare, contrast, and evaluate application composition scenarios with ADIOS, XPMEM, and TCASM

# Source Code Availability

- Git Source Code Repository:
  - git clone https://software.sandia.gov/git/nvl
  - See README for build instructions
- Build Appliance
  - Includes software dependencies needed for NVL development
  - Includes NVL checkout and pre-built NVL images
  - http://software-login.sandia.gov/~ktpedre/hobbes/
    hobbes_build_appliance.vmwarevm.tar.gz

# XPRESS: LXK/RIOS Research Goals

- XPRESS aims to increase synergy of compute node OS kernel and user-level runtime systems
  - Today: Runtime must work around host OS, assume worst case
  - Vision: Runtime cooperates with host OS, delegated more control
- Key RIOS drivers (Runtime Interface to the OS)
  - Runtime needs guarantees about resource ownership and behavior
  - OS needs way to shift resources between multiple runtimes
  - Two-way interfaces needed for key resources
    - Runtime tells OS what it needs, OS tells runtime what it gets
    - OS remembers original request, notifies runtime if more resources become available. Notifies runtime of resources need to be reclaimed.
  - Event-based protocol to notify of dynamic events (e.g., power state change, transient error)
- LXK = Kitten + RIOS

# XPRESS: Areas Covered by RIOS

- Legacy support services
- Job management
- Memory management
- Thread management
- Network interface
- System topology and locality
- Introspection
- File I/O
- Power management

**SANDIA REPORT**

SAND2013-XXXX
Unlimited Release
Printed September 2013

**Runtime Interface to the Operating System (RIOS) Specification**

Brian Barrett, Ron Brightwell, Stephen Olivier, Kevin Pedretti, Dylan Stark, and Thomas Sterling

Prepared by
Sandia National Laboratories
Albuquerque, New Mexico 87185 and Livermore, California 94550

Sandia National Laboratories is a multi-program 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.

Approved for public release; further dissemination unlimited.

# Bull eXascale Interconnect (BXI)

- New, next-generation interconnect fabric

- Hardware-based network interface offload

- Several virtual networks for QoS

- Adaptive routing

- End-to-End and link-level retry for reliability

- Based on Portals 4 communication library

http://xstack.sandia.gov/hobbes
http://xstack.sandia.gov/xpress