

HPC Operating System Research Areas and Challenges



PRESENTED BY

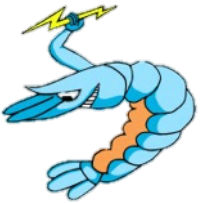
Kevin Pedretti, Ron Brightwell, Jack Lange, Andrew Younge

ASCR Roundtable Discussion on OS Research
January 25, 2021

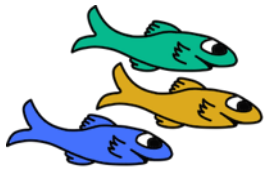
OS/R Research Approach: Build Real Systems



github.com/hobbesosr/kitten



www.prognosticlab.org
www.v3vee.org

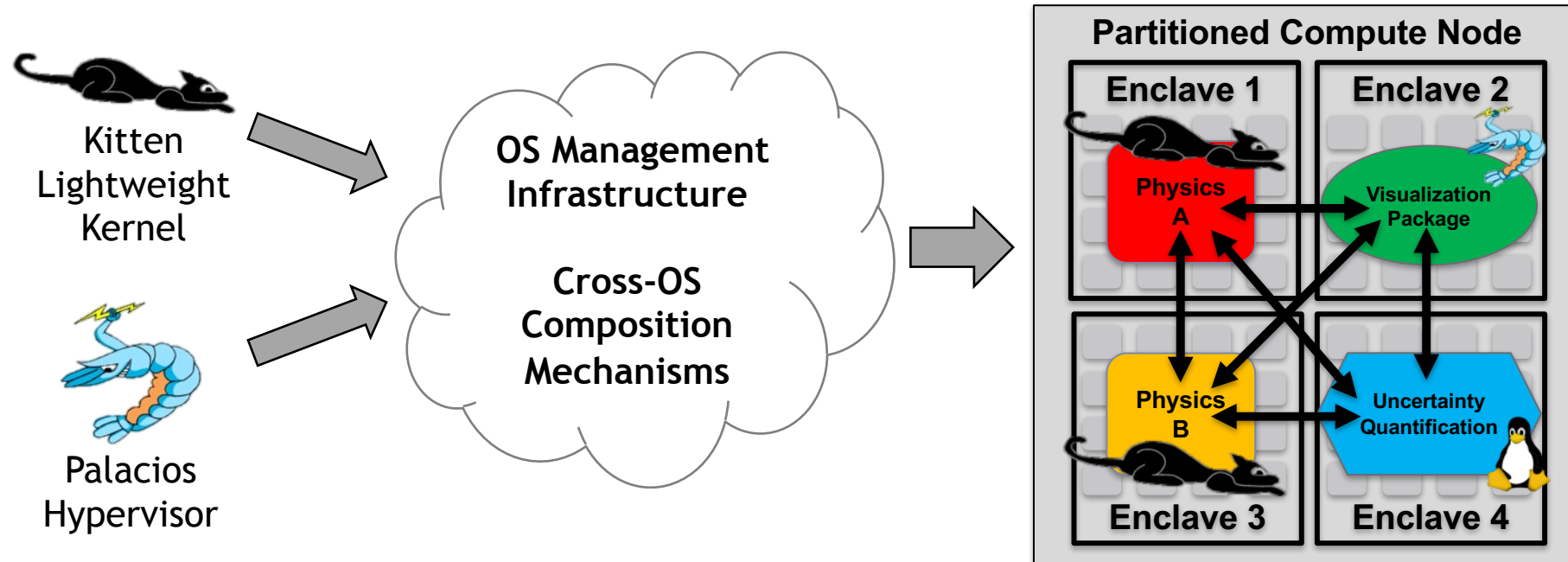


<https://github.com/HobbesOSR/nvl>

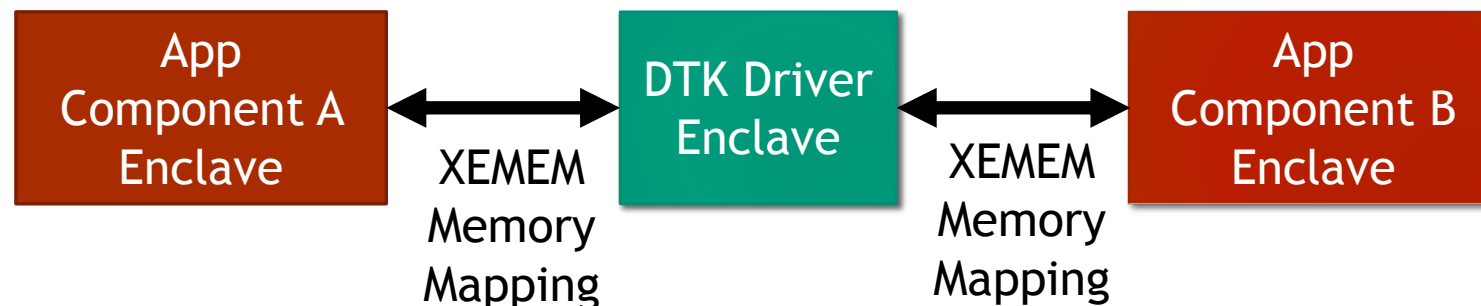
- **Kitten Lightweight Kernel**
 - SUNMOS (1993), Cougar (1997), Catamount (2004), Kitten (2008-)
 - Linux ABI + API compatible user space, compile on Linux run on Kitten
- **Palacios Virtual Machine Monitor**
 - OS independent, easily embeddable design, lightweight resource management
 - Demonstrated < 5% overhead on 4K Red Storm nodes, Kitten+Palacios (VEE'11)
- **SMARTMAP / XPMEM / XEMEM**
 - Enables processes running on same node to directly access each others memory
 - SMARTMAP for Catamount (SC'08), Cray XPMEM for Linux, XEMEM (HPDC'15)
- **Hobbes Node Virtualization Layer (ASCR Exascale OS/R FY13-17)**
 - Enables partitioning compute node resources among multiple OS/R stacks
 - Provides cross OS/R stack composition mechanisms; virt on Cray XC (Cluster'17)
 - Provides performance isolation at hardware and system software levels (HPDC'15)

Hobbes Focused on Application Component Composition

Still a Major Challenge, esp. with Extreme Heterogeneity



Key challenge is sharing + transforming data efficiently between discrete components, Approach prototyped by Hobbes using Data Transfer Toolkit (DTK):



Observations from Fielding an ARM Based Supercomputer



- After Hobbes, focused on developing Astra:



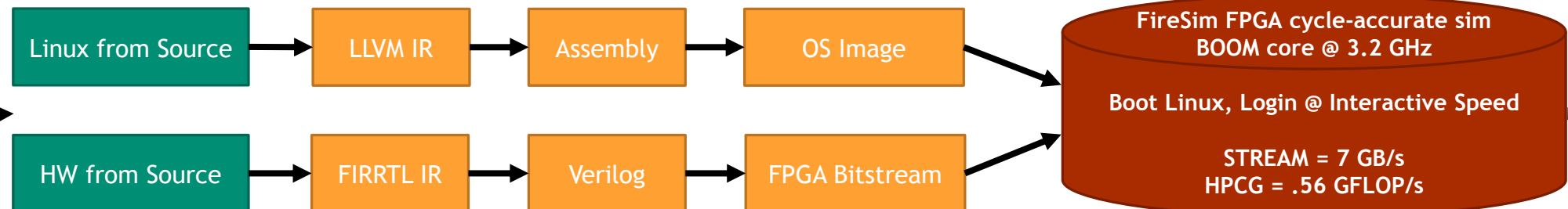
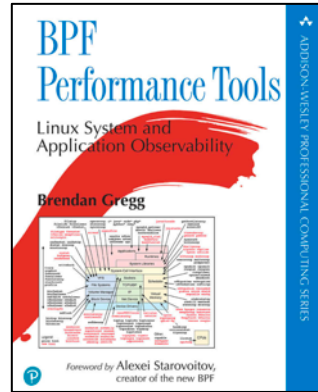
5,184 Marvell ThunderX2 28-core ARM Processors (145,152 cores)

SC'20 – [Chronicles of Astra: Challenges and Lessons](#)

- Goal to mature Arm HW+SW ecosystems for DOE / NNSA HPC workloads
- Lack of observability tools slowed us down (e.g., debug Linux, thermal issues)
- Even for CPU-only system, many user challenges with NUMA and affinity issues
- Users very interested in containers, challenging when user laptops != ARM
- System usage model predominantly ~1990's era batch scheduled SPMD model

OS/R Challenges and Research Areas

1. Application component orchestration systems for HPC workflows
2. Observability, provide users + operators with actionable info
 - BPF changing the game, C -> LLVM -> eBPF program injected into Linux Kernel
 - Observability is a pre-requisite for autonomous resource management
3. Integration of HPC, Data Intensive, and Cloud
 - Container orchestration for **long-lived services** + HPC apps; Move beyond pure batch model
 - Form collaborations with hyperscalers to design future HPC cloud technologies
4. Co-design of Extremely Heterogeneous Hardware and OS
 - Develop HW+SW interfaces to enable portable system software, reduce user burden
 - Open-source HW enabling rapid end-to-end co-design, Berkeley FireSim example:



Design Iteration in Minutes / Hours